The background of the slide is a dense, textured layer of light brown wood chips or mulch, with individual chips of varying lengths and orientations creating a complex, fibrous pattern.

Genomic and Bioinformatic Technologies

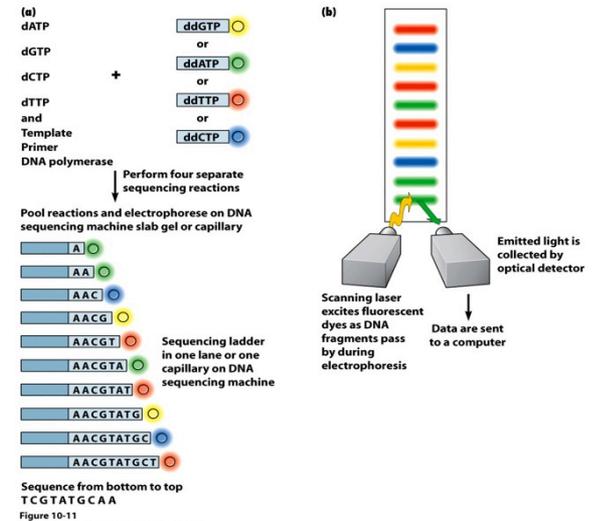
Jeremy Schmutz, HudsonAlpha Institute
Todd Mockler, Donald Danforth Plant Science Center
DOE ARPA-E “Advanced Plant Phenotyping for
Increased Bioenergy Crop Yields”

Overview

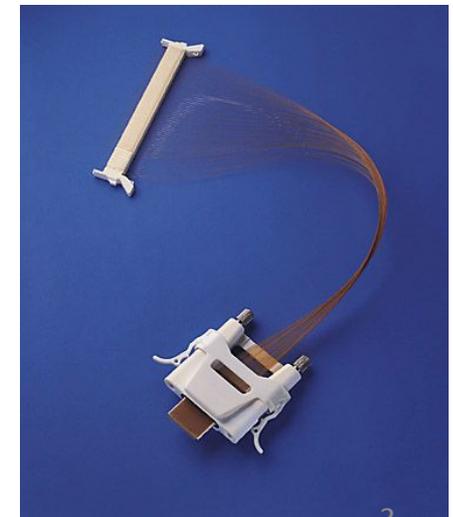
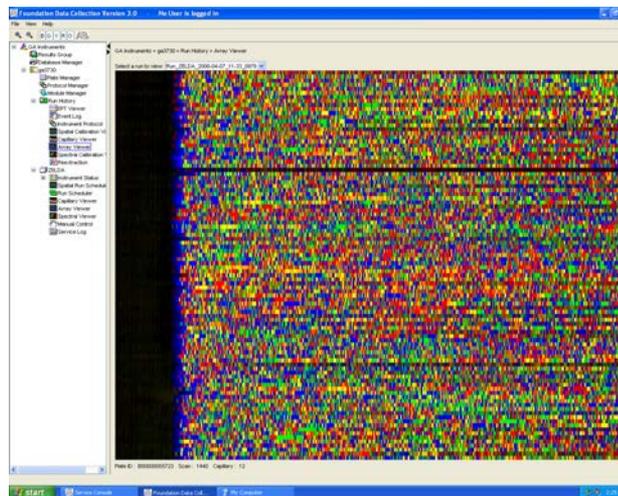
- Sequencing technology changes
- Plant reference genomes
- Sequence based tools for interrogation of reference genomes
- Plant databases/computational analysis platforms
- Bioinformatics landscape
- Controlled environment phenotyping and analytics
- Where we need to go

Sanger sequencing (3730xl 2002)

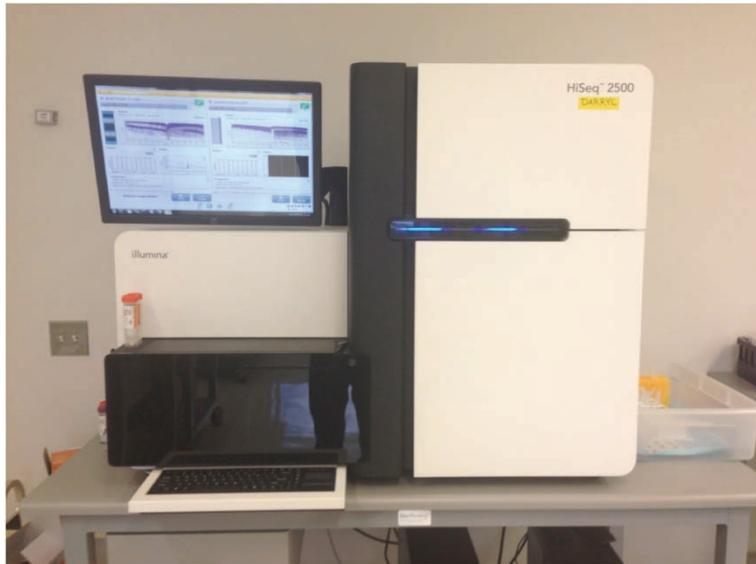
- Integrated robot loader and plate piercer –holds 24 plates
- 2 hour run time per 96 samples, 800 bp reads@ ~900k bps per day – ran unattended
- \$937,500 per GB



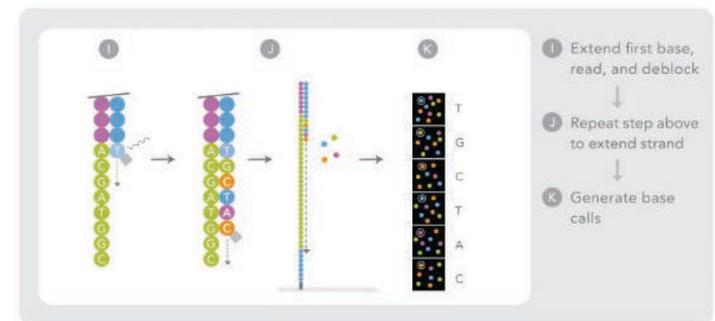
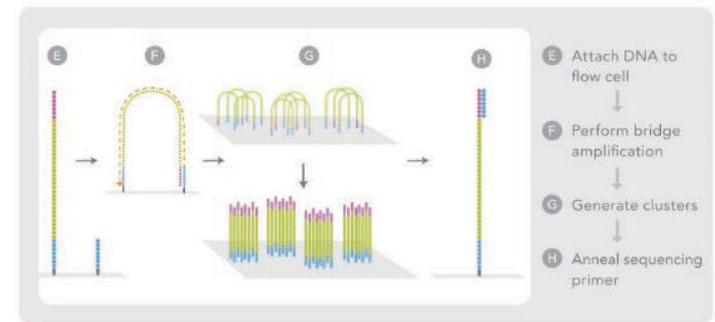
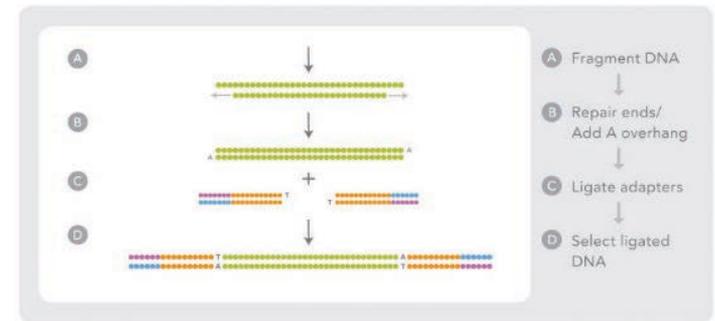
Recombinant DNA: Genes and Genomes (3e)
© 2007 W. H. Freeman and Company



Modern data collection



- Produces ~900 Gb of data in 6 days (150 Gb per day)
- ~40k per run, or \$44 per GB excluding library, machine cost, and data handling (1,500x Sanger & 23,000x lower cost)
- 2x125 reads currently, can also generate 2x200 reads at 2x the cost and ¼ the throughput



©2008, Illumina Inc. All rights reserved

@97ZXTR1:359:H9E8NADXX:1:1101:5045:3839 1:N:0:CGATGT
GCCCTGCCTCCGAGCCAGATTCCTTCCCGTTCGTTCCCGTGGCTTCTGGAGTCTTCTAGATGATAGAAA
TTGCGCGACACGTTAATATCTCTATGTAAACCCGACGTGTGGGCCTTT

The background of the image is a close-up, vertical view of weathered wood. The wood is heavily textured, showing deep grooves, splintered edges, and a mix of dark brown and greyish-blue tones. The layers of wood are clearly visible, creating a complex, fibrous pattern. A semi-transparent, light-colored rectangular box is centered horizontally across the middle of the image, containing the text "Plant reference genomes" in a bold, black, sans-serif font.

Plant reference genomes

Why *de novo* plant genomes are hard

Polyploidy & Polymorphism

Repeat content & Transposons



**MECHANISMS THAT DRIVE
PLANT EVOLUTION!**



33%

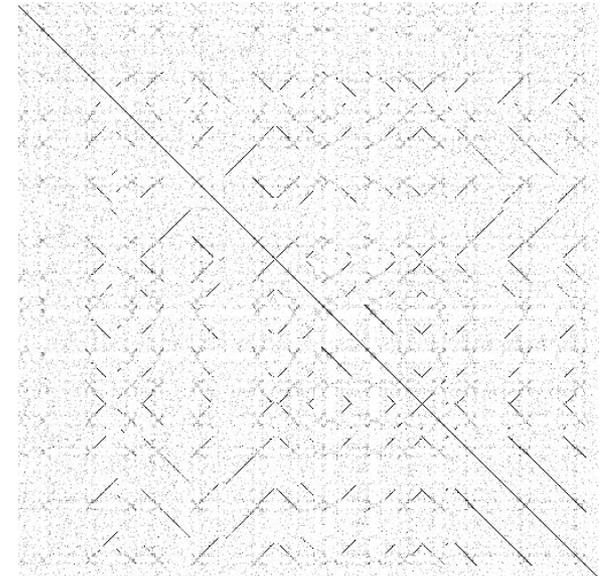
34%

44%

46%

47%

GC

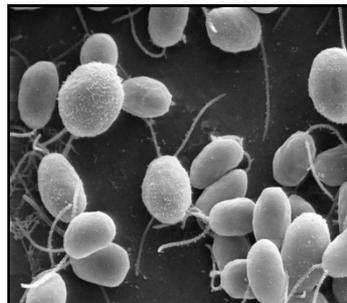


Reference plant genomes

- Come in different “flavors” – quality and completeness (finished/ improved vs. draft vs. drafty draft)



Sanger based BAC by BAC projects



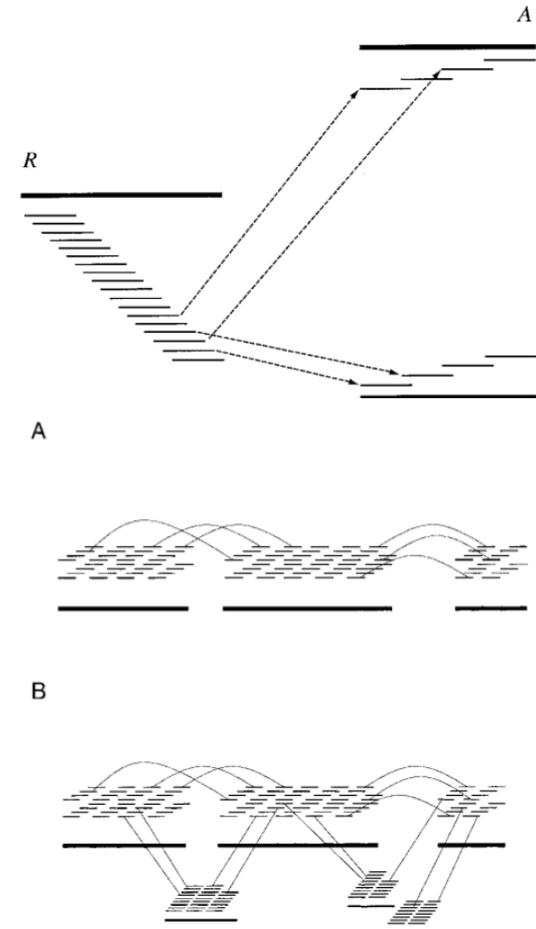
Sanger based WGS genomes



NGS references

WGS references

- Sequence the whole genome at once using multiple sized libraries
- Calculate all overlaps between reads and then work out a parsimonious path through the data set
- **WGS benefits**
 - Small number of good subclone libraries needed
 - Short Subclone (3-4kb)
 - Medium Subclone (6-9kb)
 - Long Fosmid (32-36kb)
 - Longer BAC (100-160kb)
 - No mapping phase
 - Less complex pipeline
- **WGS drawbacks**
 - Inconsistent genome coverage
 - Difficult to assess quality
 - Laborious post-production phase
 - Don't know what you have till you are done



ARACHNE: A Whole-Genome Shotgun Assembler

Serafim Batzoglou,^{1,2,3} David B. Jaffe,^{2,3,4} Ken Stanley,² Jonathan Butler,² Sante Gnerre,² Evan Mauceli,² Bonnie Berger,^{1,5} Jill P. Mesirov,² and Eric S. Lander^{2,6,7}

NGS “reference” genomes

- Sequence fragments of short inserts of varying sizes, recombination pairs, collapse very large data set into consensus

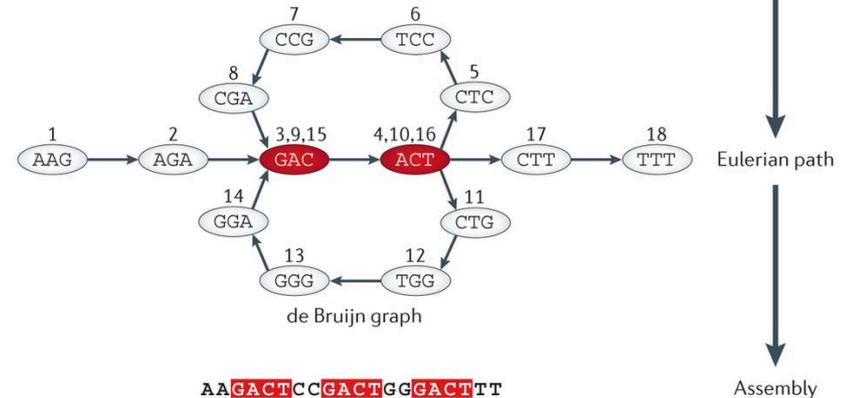


- **NGS benefits**

- No cloning
- Data collection is inexpensive
- Sequence is high quality

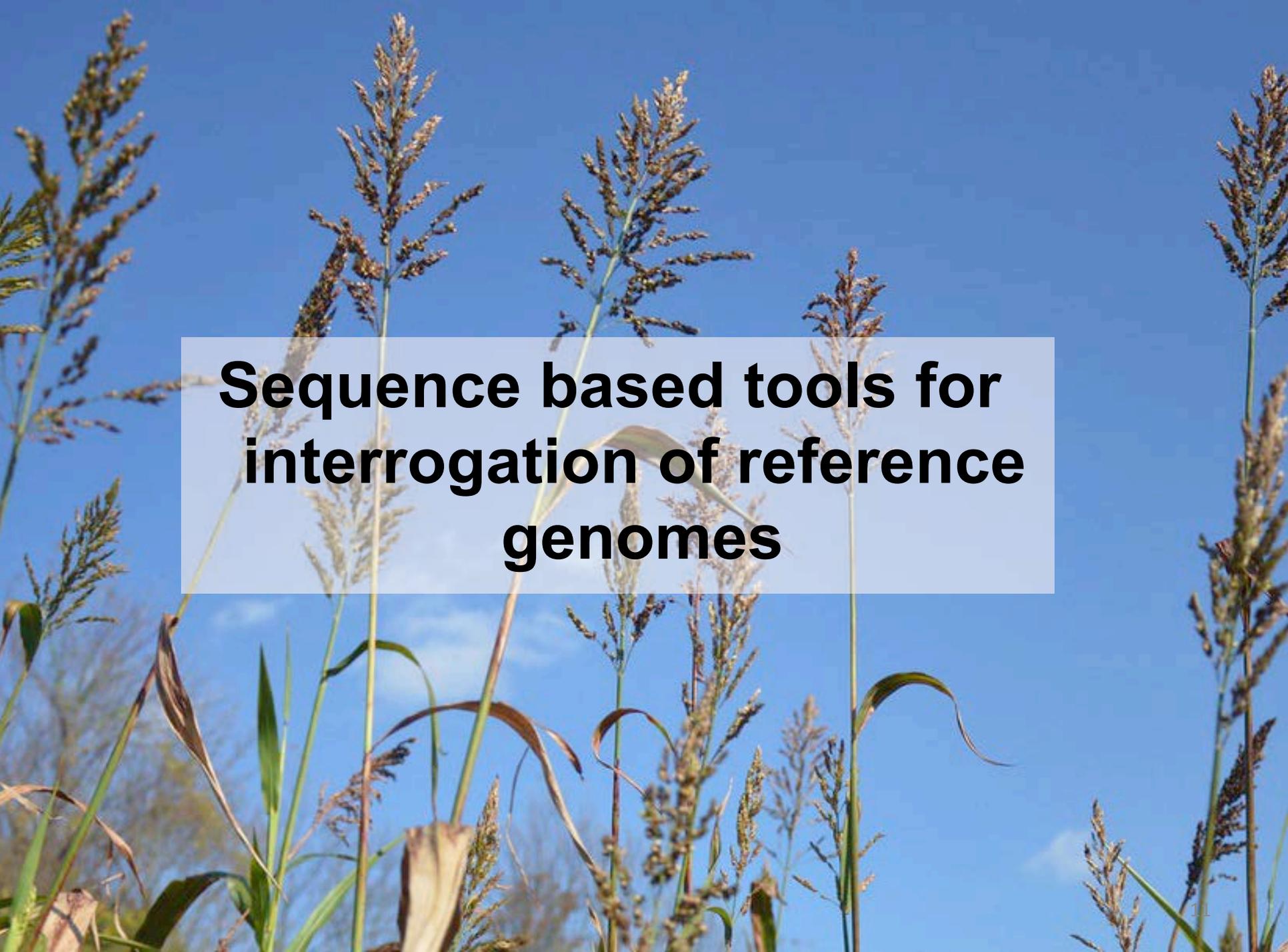
- **NGS drawbacks**

- Biases with whatever technology is used
- Difficult to assess quality/completeness
- Tend not to include long pairs, mapping information, validation
- Lots of contigs, low repeat content, complex to use in practice
- Difficult to improve with directed means
- Fragmented



AA**GACT**CC**GACT**GG**GACT**TT

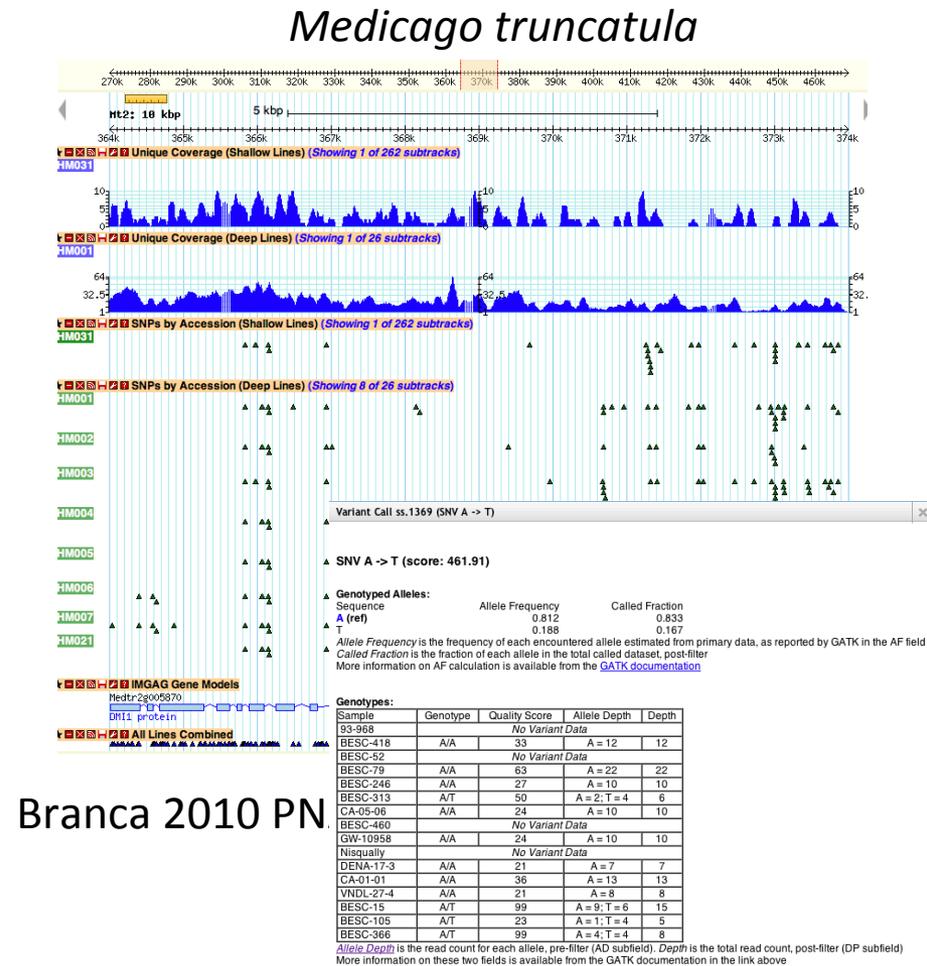
Computational solutions for omics data
Bonnie Berger, Jian Peng & Mona Singh
Nature Reviews Genetics 14, 333–346 (2013) doi:
10.1038/nrg3433

A photograph of a field of tall grasses, likely sorghum, with their seed heads reaching towards a clear blue sky. The grasses are in the foreground and middle ground, creating a sense of depth. The text is overlaid on a semi-transparent white box in the center of the image.

**Sequence based tools for
interrogation of reference
genomes**

Resequencing for SNP/indel identification

- Collect DNA from cultivars, breeding populations, make frag libraries, resequence
- 15 x of inbred sorghum = 14 Gb, \$1,000 (detect homozygous changes)
- 30x of outbred sorghum = \$1,800 (detect heterozygous changes)
- Lower coverage for pedigree typing 2x = \$200



Branca 2010 PN

SnpEff Variant Annotation, where available:

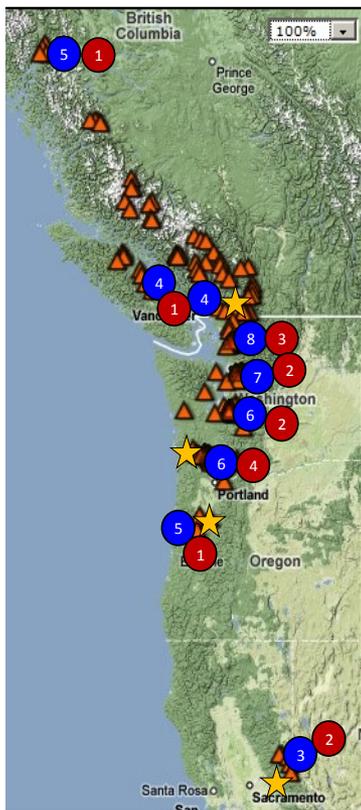
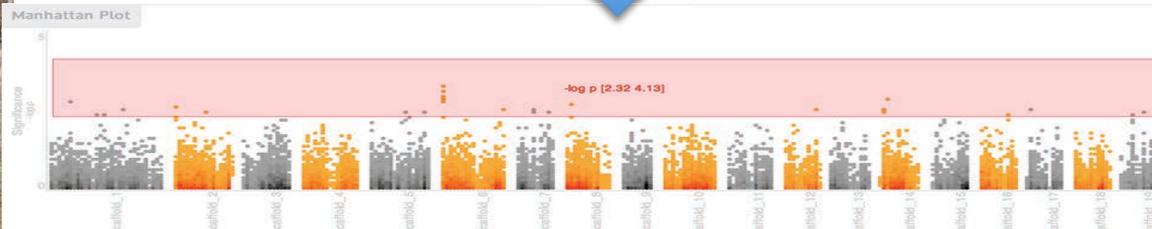
Effect	Effect Impact	Codon Change	Gene Name	Transcript	Exon	Genotype
DOWNSTREAM	MODIFIER	1516	Potr.001G001900	Potr.001G001900.1		T
DOWNSTREAM	MODIFIER	2857	Potr.001G002000	Potr.001G002000.1		T
DOWNSTREAM	MODIFIER	2857	Potr.001G002000	Potr.001G002000.2		T
INTRON	MODIFIER		Potr.001G001800	Potr.001G001800.1	1	T
INTRON	MODIFIER		Potr.001G001800	Potr.001G001800.2	1	T
INTRON	MODIFIER		Potr.001G001800	Potr.001G001800.3	1	T
UPSTREAM	MODIFIER	1535	Potr.001G001700	Potr.001G001700.1		T

Real world example

Dense SNPs
across
population

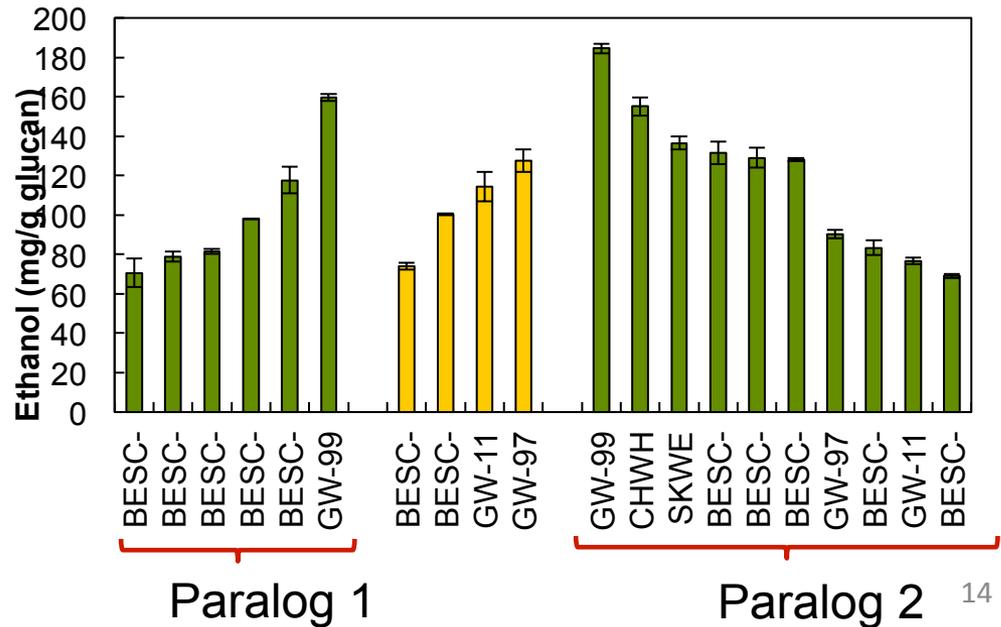
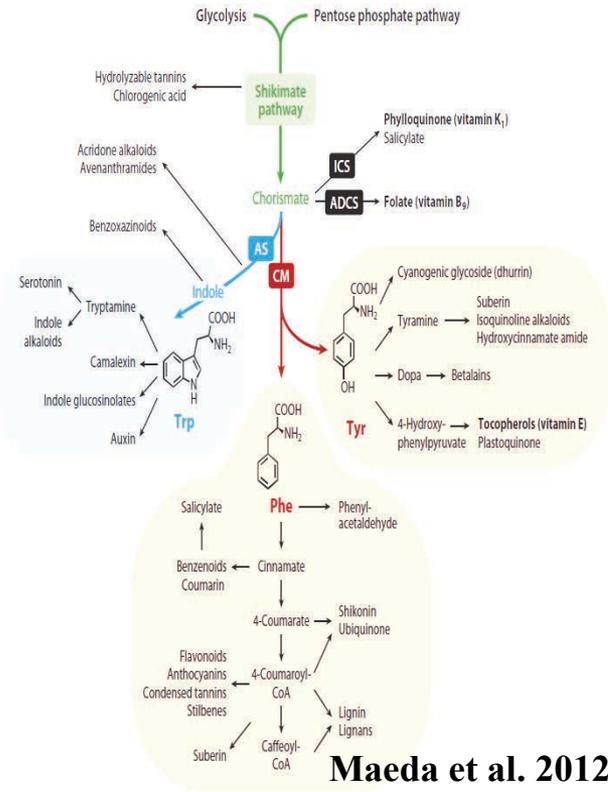
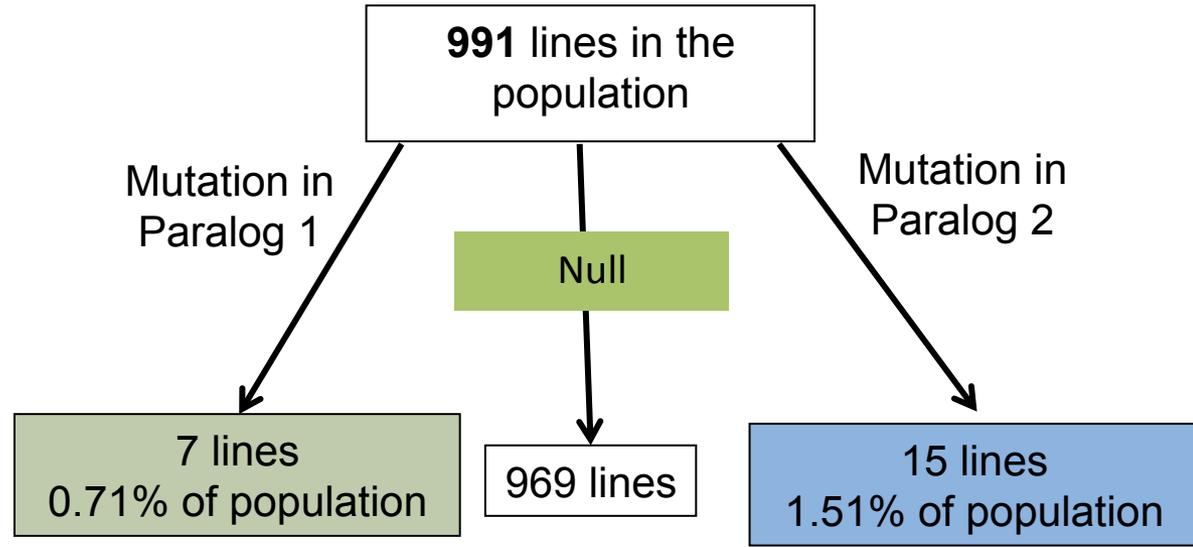
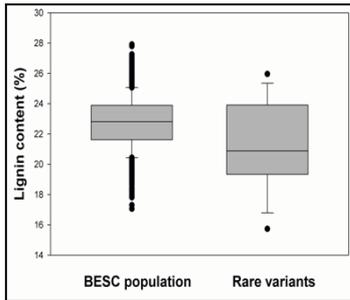
#CHROM	POS	REF	201782_400122	201782_400194	
	201782__400495	...			
Chr01	4	C	C/C	C/C	C/T
Chr01	5	A	A/A	A/A	A/A
Chr01	7	ACCCC	ACCCC/ACCCC	ACCCC/ACCCC	ACCCC/ACC
Chr01	18	A	A/A	A/A	A/A
Chr01	20	A	A/A	A/A	A/A
Chr01	21	C	C/C	C/C	C/C
Chr01	22	C	C/C	C/C	C/C
Chr01	23	C	C/C	C/C	C/C
Chr01	24	CAAA	CAAA/CAAA	CAAA/CAAA	CAAA/CAAA
Chr01	34	ACCCC	ACCCC/ACCCC	ACCCC/ACCCC	ACCCC/ACCCC
....					

Detailed
phenotyping



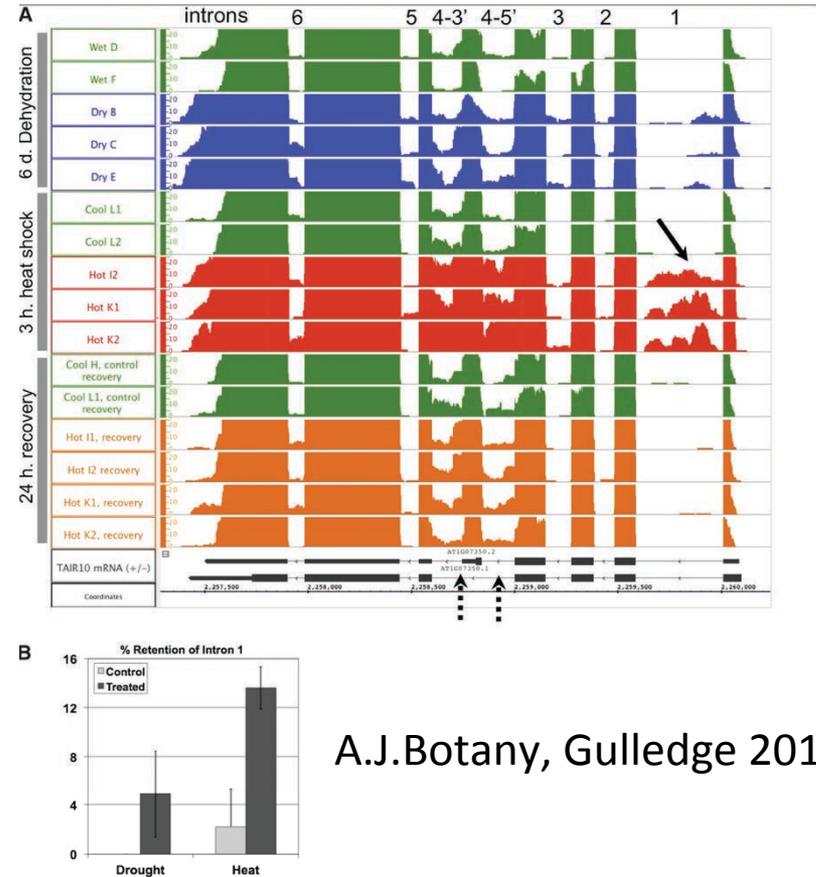
Jerry Tuskan: ORNL/BESC

Lignin quality traits



RNA-seq

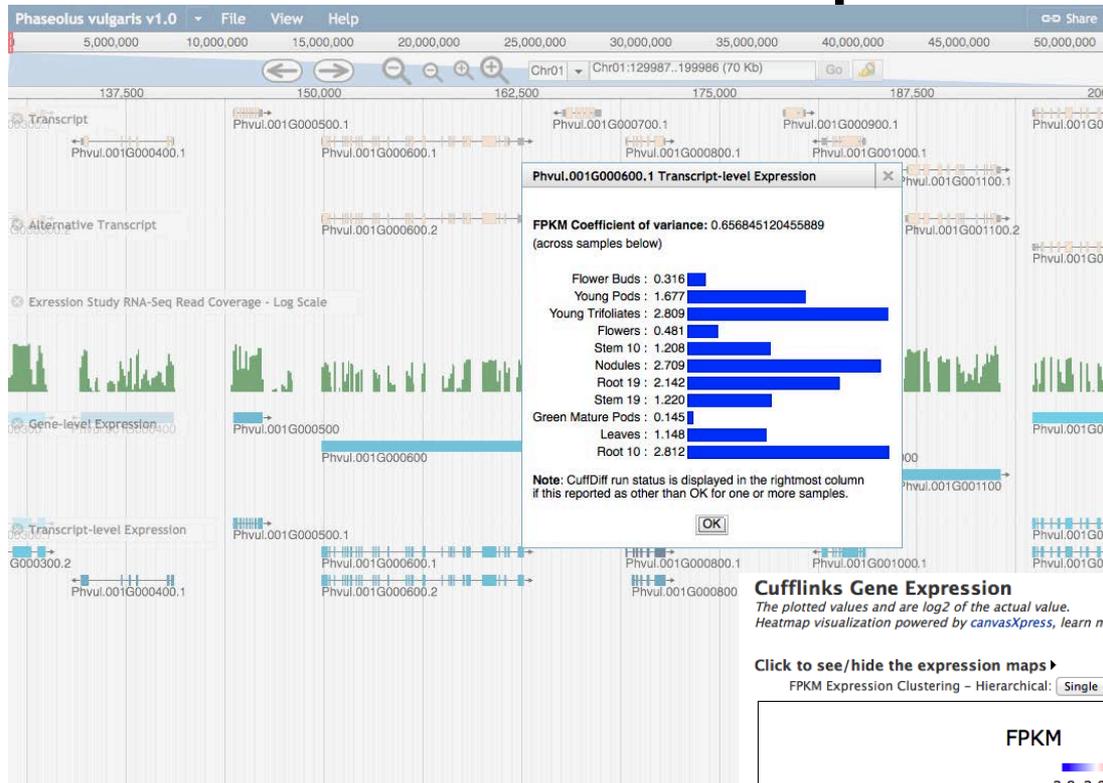
- Isolate RNA, make RNA-seq (RT) libs, collect data
- Expression analysis, tissue to tissue, condition to condition, or genotype to genotype comparisons
- Typically for full analysis run at 10-20M pairs per (or 10-20 samples per lane ~\$325 per samples
- Variations such as 3' counting, non-coding RNA collection, microRNA collection



A.J.Botany, Gulledge 2012

Arabidopsis: alternative splicing in splicing regulator SR45a

RNA-seq data displays



Cufflinks Gene Expression

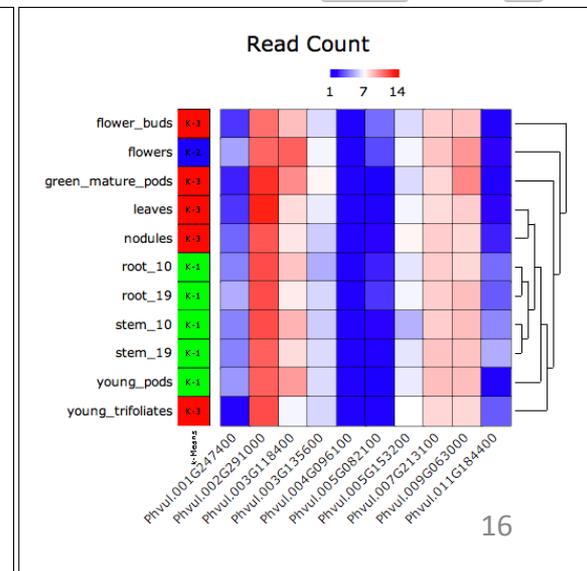
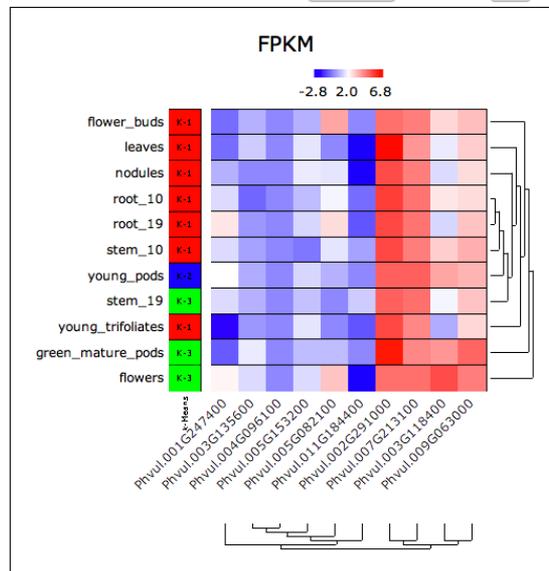
The plotted values are \log_2 of the actual value.
Heatmap visualization powered by *canvasXpress*, learn more about the *display options*.

Click to see/hide the expression maps ▶

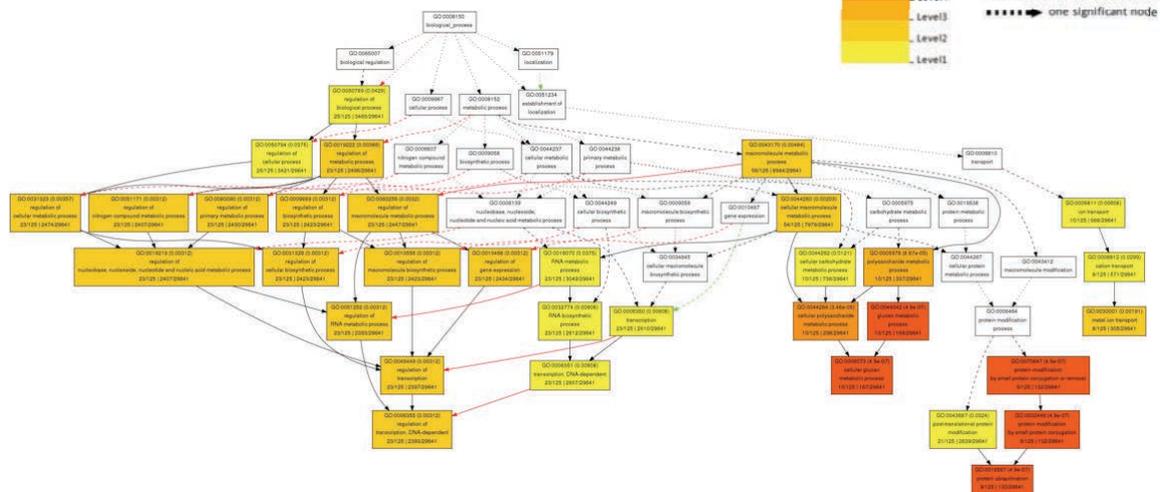
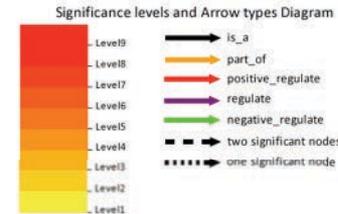
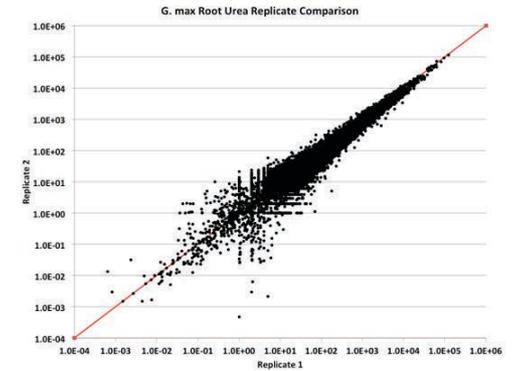
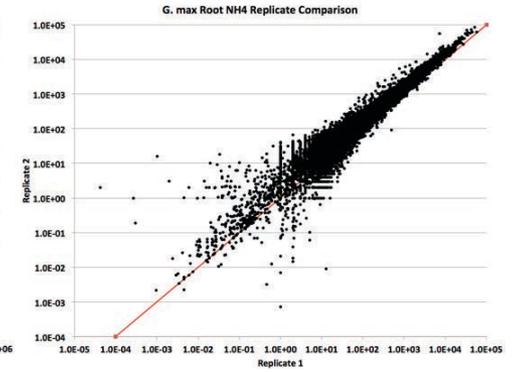
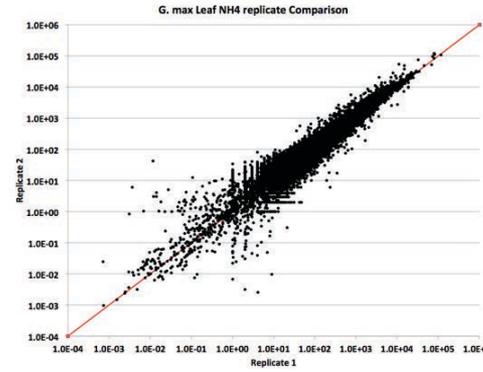
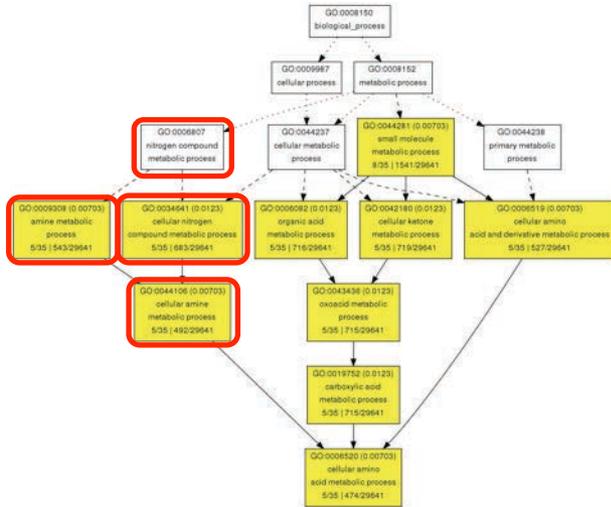
FPKM Expression Clustering - Hierarchical: Single and K-means: 3

Count Expression Clustering - Hierarchical: Single and K-means: 3

Phaseolus vulgaris

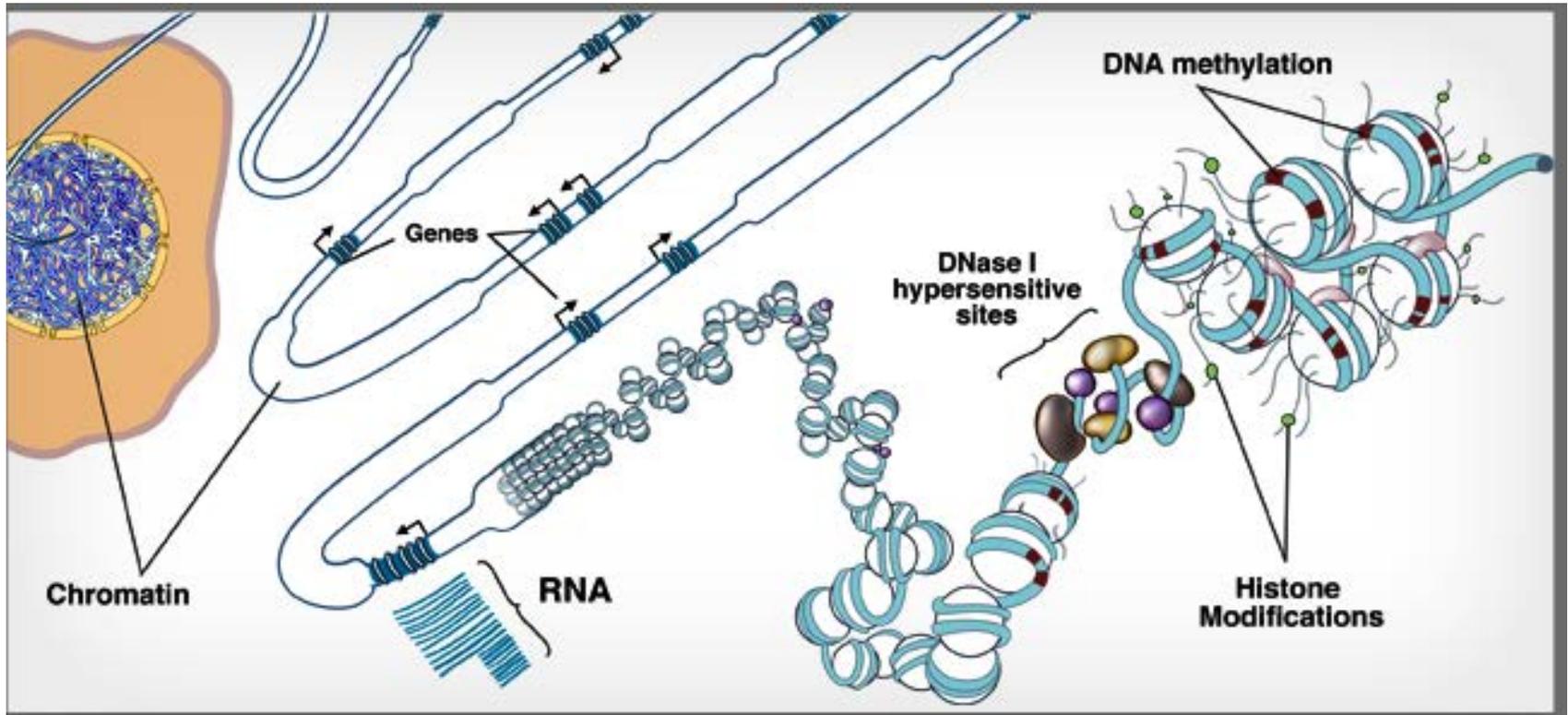


RNA-seq for expression comparisons

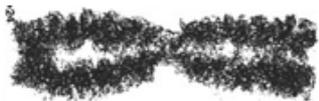


Glycine max (soybean)
w/ Gary Stacey : Mizzou

Histone modifications



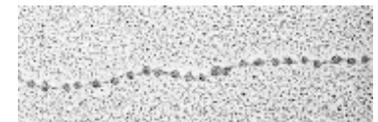
<http://www.roadmapepigenomics.org>



Long range fiber-fiber interaction



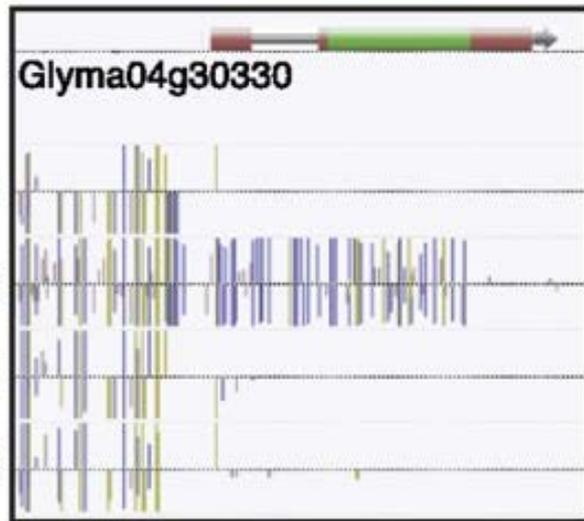
Short range internucleosomal interactions



Nucleosome

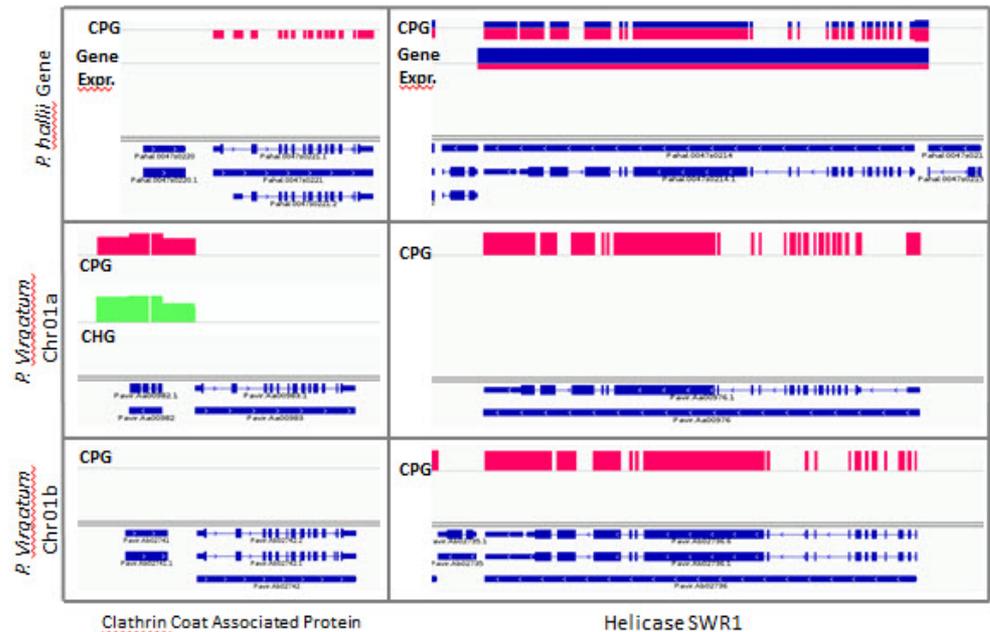
5mC methylation is likely important for crops!

- Bisulfate sequencing for bp resolution or Reduced Methylation Assay for broad picture
- Inherited methyl states affect transcription



Schmitz, GR, 2013

Glycine max (soybean)



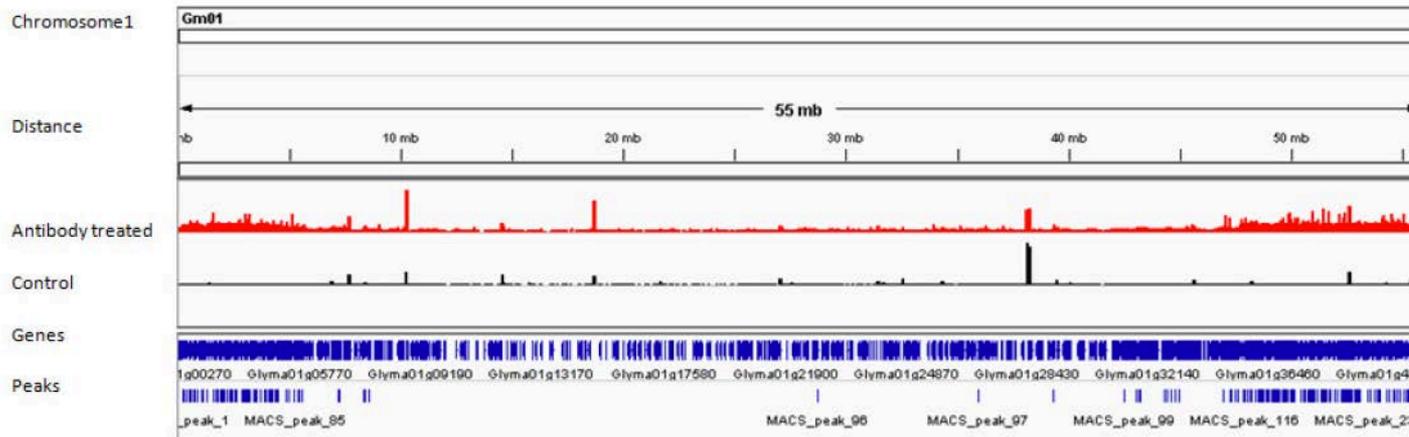
Methylation comp. between diploid *Panicum hallii* and tetraploid *Panicum virgatum*

Tom Juenger : UT-Austin

ChIP-Seq in plants

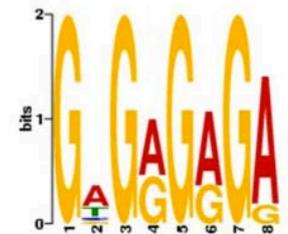
- Use antibodies to identify positions in a genome where a specific transcription factor binds

Glycine max (soybean) :NAC transcription factor

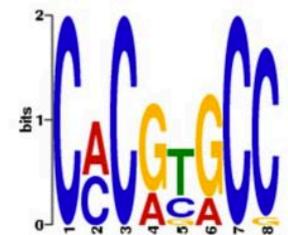


Shamimuzzaman, 2013 BMC Genomics

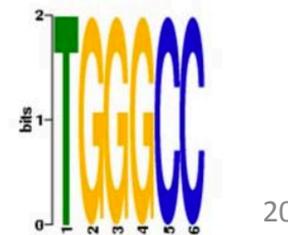
(a) Potential NAC TF Binding Motif1



(b) Potential NAC TF Binding Motif2

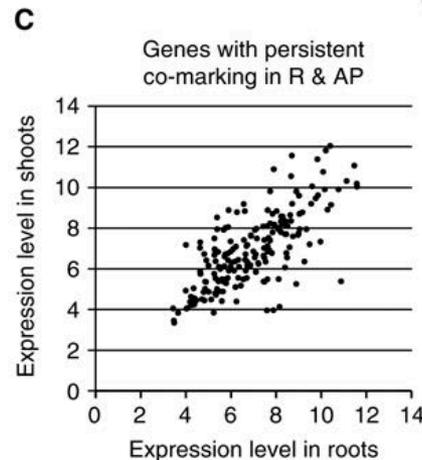
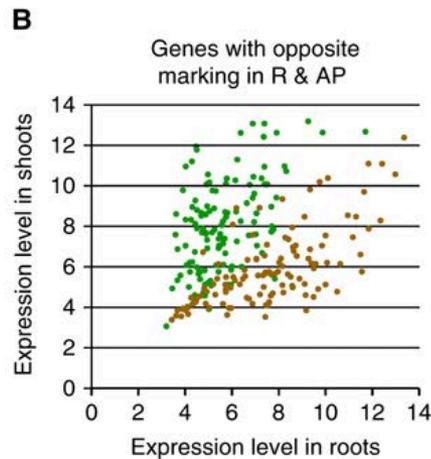


(c) Potential NAC TF Binding Motif3

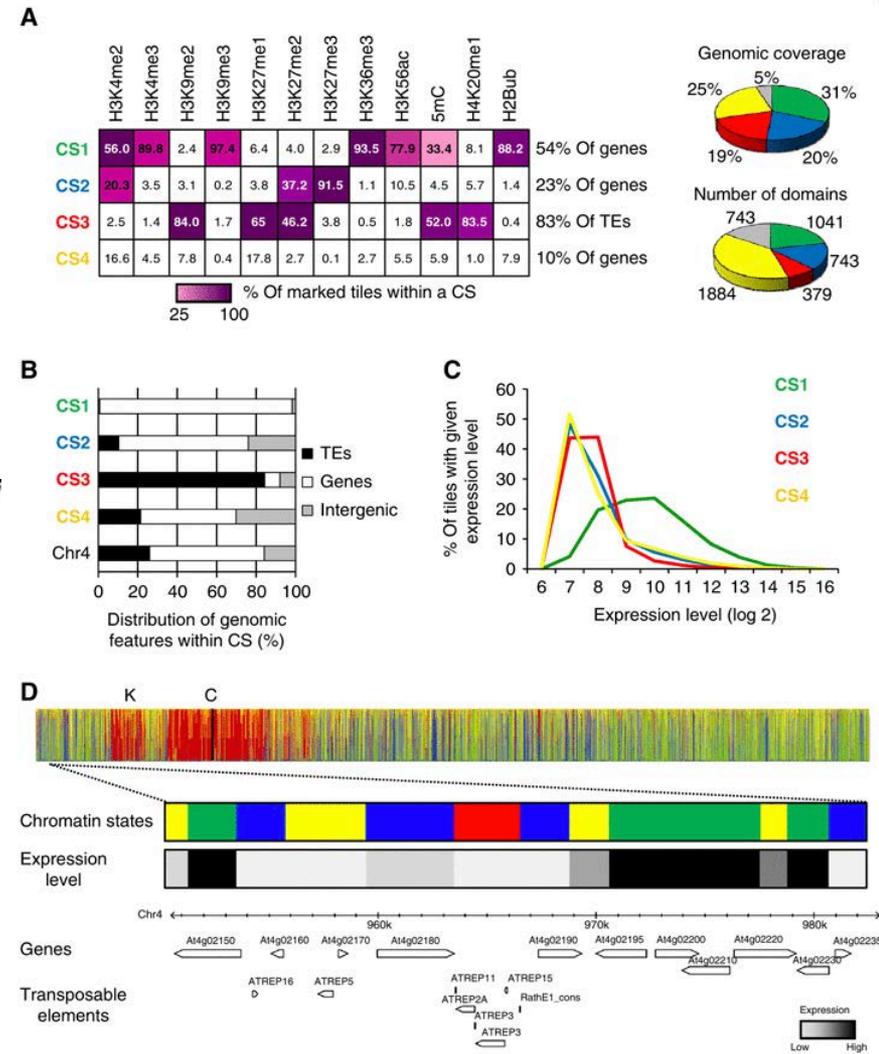


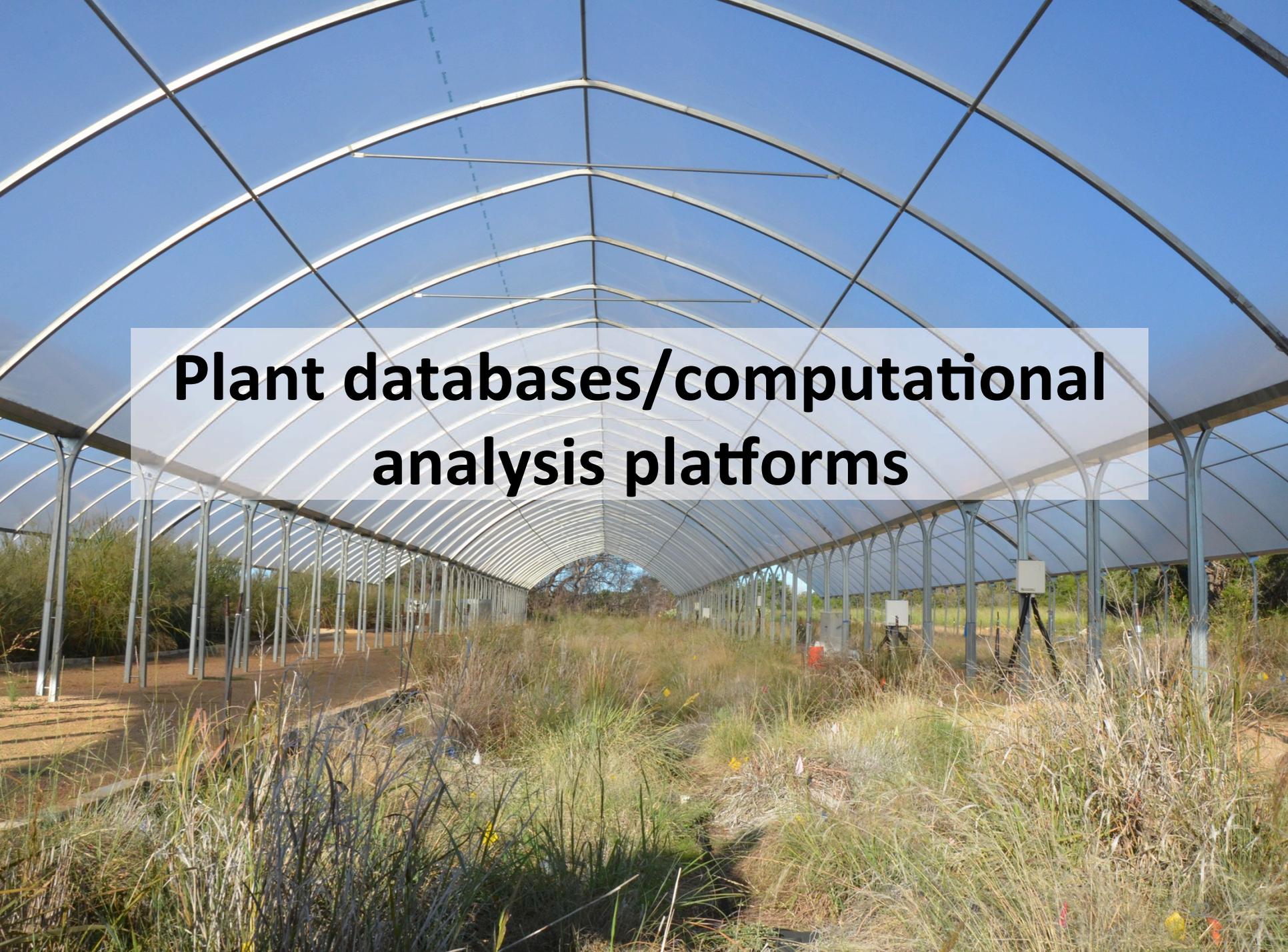
Histone modifications in plants

- Identify histone modifications using ChIP-Seq
- Index active genes, repressed genes, silent repeat elements and intergenic regions.



Arabidopsis thaliana





**Plant databases/computational
analysis platforms**

JGI: Phytozome Portal

Welcome to phytozome The JGI Comparative Plant Genomics Portal

Phytozome quick start

Explore a JGI flagship genome



Glycine max Wm82 a2.v1 Setaria italica v2.1 Populus trichocarpa v3.0 Physcomitrella patens v3.0 Chlamydomonas reinhardtii v5.5 Brachypodium distachyon v1.2 Panicum virgatum v1.1 Sorghum bicolor v2.1

Query: Enter keywords or sequence

Select from all species/nodes in Phytozome

Early release species

Help with Phytozome

Video tutorials

- Video tutorials will be available shortly

About Phytozome

Phytozome is the Plant Comparative Genomics portal of the Department of Energy's Joint Genome Institute. Families of related genes representing the modern descendants of ancestral genes are constructed at key phylogenetic nodes. These families allow easy access to clade-specific orthology/paralogy relationships as well as insights into clade-specific novelties and expansions. As of release v10, Phytozome provides access to forty-six sequenced and annotated green plant genomes which have been clustered into gene families at 13 evolutionarily significant nodes. Each gene has been annotated with PFAM, KOG, KEGG, PANTHER and GO assignments, where possible. Query-based data access is provided by Phytozome's InterMine and BioMart instances, while bulk data sets can be accessed via the JGI's Genome Portal (registration required). JBrowse genome browsers are available for all genomes.

News (details...)

(2014-04-27) v10.0.1 bug fix release

(2014-04-04) Brachypodium v2.1, Arabidopsis halleri v1.1, Panicum hallii v0.5 early releases

(2014-03-18) Phytozome v10.0 beta is now live!

System Status (2014-05-28 11:36)

- Search
- BLAST
- BLAT
- InterMine
- Database

Data and Analyses available in v10

- 46 annotated plant and algal genomes
- Protein families constructed at
- InterPro/KEGG/KOG proteome annotation
- natural diversity data for *Ptr*, *Egr*, *Bdi*
- expression data for *Cre*, *Pvu*, *Gma*

Phytozome Features

- Genome-centric, gene-centric, family-centric search and visualization
- JBrowse genome browsers
- InterMine queryable data warehouse
- Searchable by sequence similarity, GO/PFAM/KEGG/KOG/PANTHER, symbol, identifier, keyword,
- Custom user lists, data downloads
- Programmatic data access via web services API and multi-language client libraries

Family chloroplastic drought-induced stress protein of 32 kD

Family Info

Identifier Poplar-Malvidae gene family 46334641, 11 members
Size 11 members
Membership Ptr Cpa Gra Tca Ath Bst Bra Cru Esa Ccl
 2 1 1 1 1 1 1 1 1 1 1
KOG Class CELLULAR PROCESSES AND SIGNALING [0] : Posttranslational modification, protein turnover, chaperones

Genes in Family	Functional Annotation	MSA	Family History
<input type="checkbox"/> M Views Org ID Alias/Symbol Define Domains Synteny Exons			
<input type="checkbox"/> F G B Ptr Potri.002G016... POPTR_0002s017...			
<input type="checkbox"/> F G B Ptr Potri.005G245... POPTR_0005s267...			
<input type="checkbox"/> F G B Cpa evm.model.su...			
<input type="checkbox"/> F G B Gra Gorai.009G19...			
<input type="checkbox"/> F G B Tca Thecci.EG034...	Chloroplastic drought-induced stress pro...		
<input type="checkbox"/> F G B Ath AT1G76080.1 ATCDS32	chloroplastic drought-induced stress pro...		
<input type="checkbox"/> F G B Bst Bostr.20129s0...			
<input type="checkbox"/> F G B Bra Brara.G03401...			
<input type="checkbox"/> F G B Cru Carubv100207...			
<input type="checkbox"/> F G B Esa Thhalv100189...			
<input type="checkbox"/> F G B Ccl Ciclev100016...			

iPlant



The iPlant Collaborative

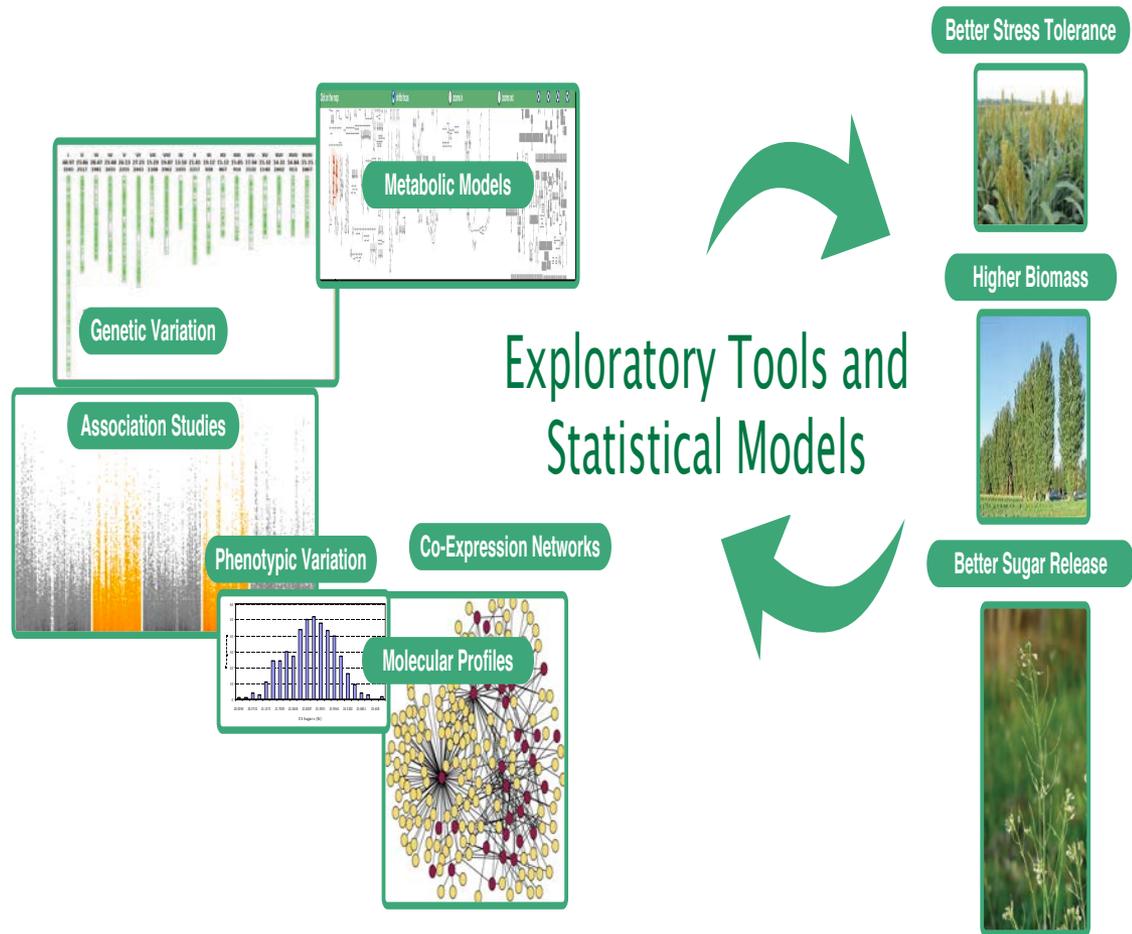
The iPlant Collaborative develops cyberinfrastructure and computational tools to solve Grand Challenges in plant science

CHALLENGE	DISCOVER	LEARN	CONNECT
<p>iPlant Genotype to Phenotype (iPG2P)</p> <p>Mapping the links between genotypes and phenotypes</p> 	<p>Discovery Environment</p> <p>Access iPlant tools through a single user-friendly interface</p> <p>MORE...</p>  <p>LOGIN</p>	<p>Upcoming Events</p> <ul style="list-style-type: none">iPlant Tools and Services Workshop, UC Davis. March 12th and 13th 2012 March 12 2012 - March 13 2012iPlant Workshop, UC Berkeley/USDAARS Albany March 12 2012 <p>MORE...</p> <p>the iPlant™ Leaflet</p> <p>Current Issue: December 05 2011</p>	<p>People at iPlant</p> <p>Community driven science</p> 
<p>iPlant Tree of Life (iPToL)</p> <p>Understanding the phylogenetic relationships between all plant life</p> 	<p>DNA Subway</p> <p>An educator-tailored interface for bringing iPlant to the classroom</p> <p>MORE...</p>  <p>LOGIN</p>	<p>News and Announcements</p> <ul style="list-style-type: none">Trellis: Climb on for Faster Access to iPlant Resources2.0 Release of Taxonomic Name Resolution Service <p>MORE...</p>	<p>My-Plant.org</p> <p>iPlant social networking</p> 
<p>Seed Projects</p> <p>Supporting diverse cyberinfrastructure needs</p> 	<p>Atmosphere</p> <p>An integrative, private, self-service cloud computing platform</p> <p>MORE...</p>  <p>LOGIN</p>		<p>Find Us...</p> 

[Register for iPlant Tools](#)

KBASE

- DOE-BER funded project to provide access to DOE computational resources
- Targeted towards metabolic modeling and statistical modeling
- Command line and GUI interfaces
- Under construction at <http://kbase.us>



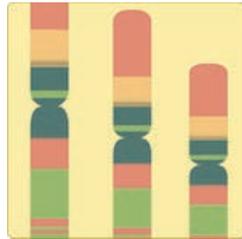
Gramene: a Comparative Resource for Plants



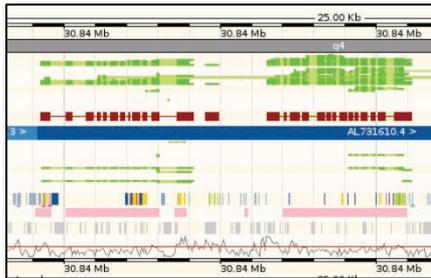
Genome Browsers

40 Species

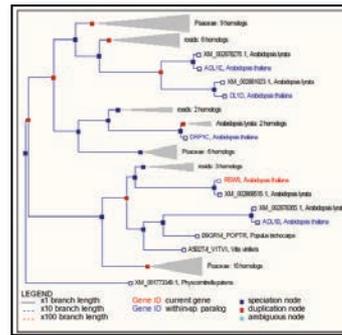
Includes Rice, Maize & Arabidopsis



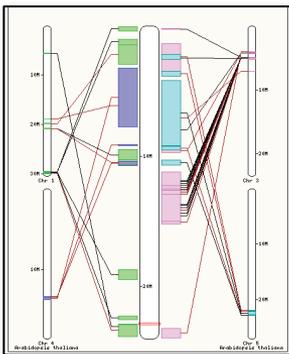
Annotation



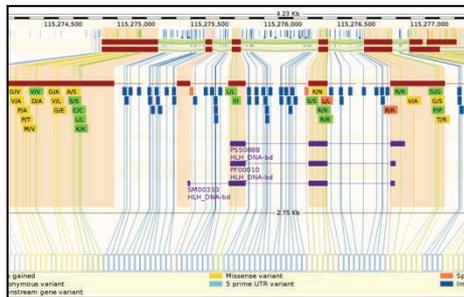
Gene trees & orthology



Comparative maps

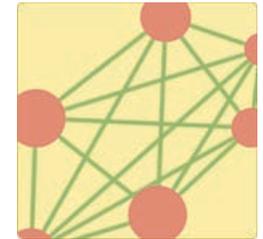


Variation & effect prediction



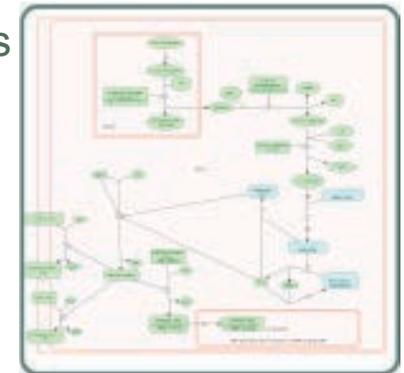
Collaborator: Paul Kersey (EBI) Ensembl

Pathways



Plant Reactome
Rice curated pathways
Maize & Arabidopsis
projections

10 species "Cyc"
metabolic pathways



Collaborator: Lincoln Stein (OICR)

Gramene services & programmatic access:
BLAST, BioMart, API, RESTFUL, MySQL

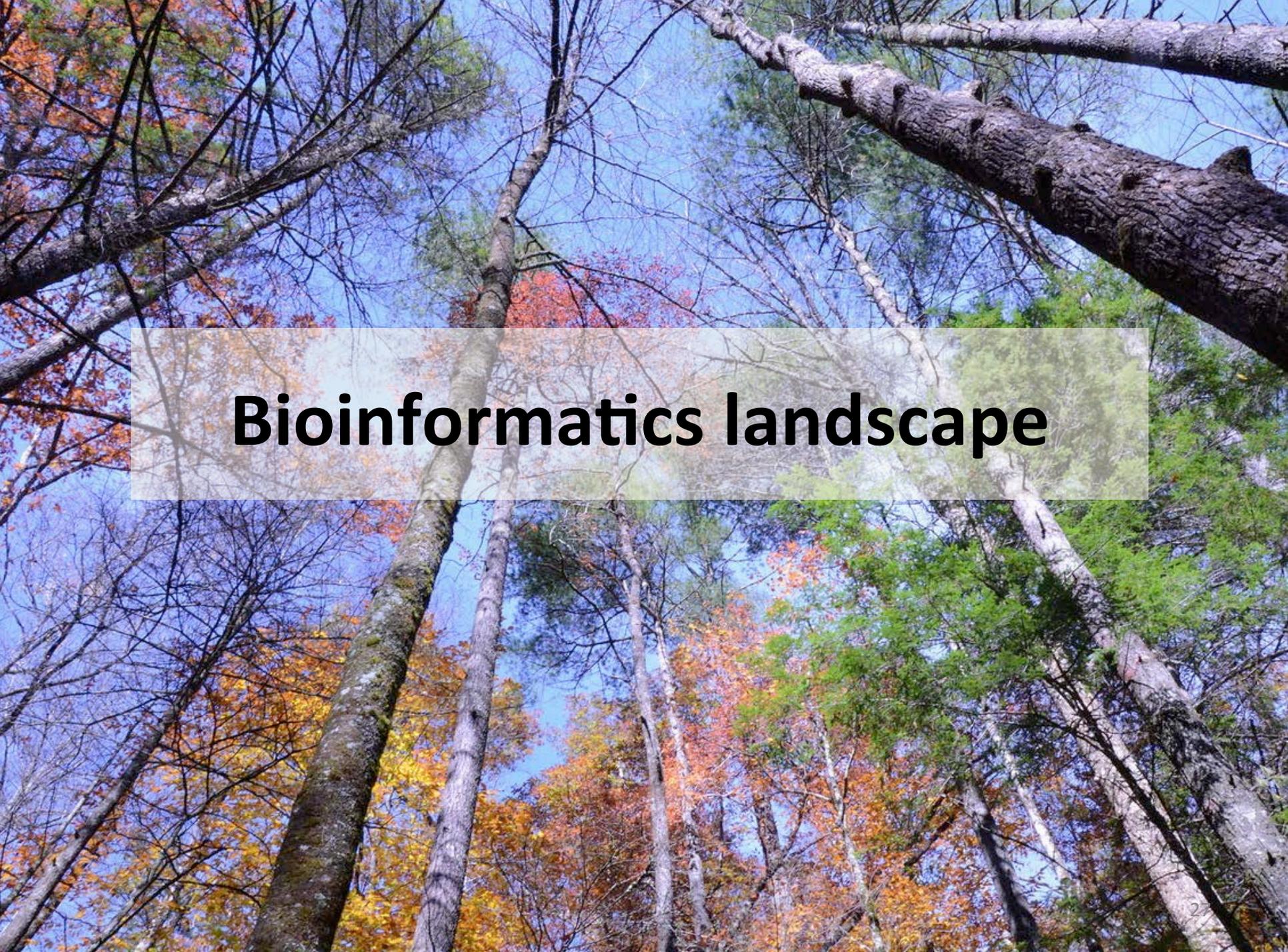
PI: Doreen Ware (CSHL)



Co-PI: Pankaj Jaiswal (Oregon State U)



Funded by NSF #0703908

A low-angle photograph of a forest. The trees are tall and thin, reaching towards a clear blue sky. The foliage is a mix of green and autumn colors like orange and yellow. The perspective is looking up from the forest floor.

Bioinformatics landscape

The Bioinformatics/Genome Analytics Landscape

Genome Assembly

- Allpaths-LG
- MaSuRCA

Gene Prediction

- MAKER-P
- Augustus

Variation Analysis

- GATK suite
- TASSEL

Network Prediction

- WGCNA
- Aracne

Transcriptome Assembly

- Trinity
- Abyss

Functional Annotation

- orthoMCL
- InterProScan

Transcriptome Profiling

- “Tuxedo” suite
- GeneCounter

Other HTS ‘Omics Analysis

- MACS
- BS-Seeker2

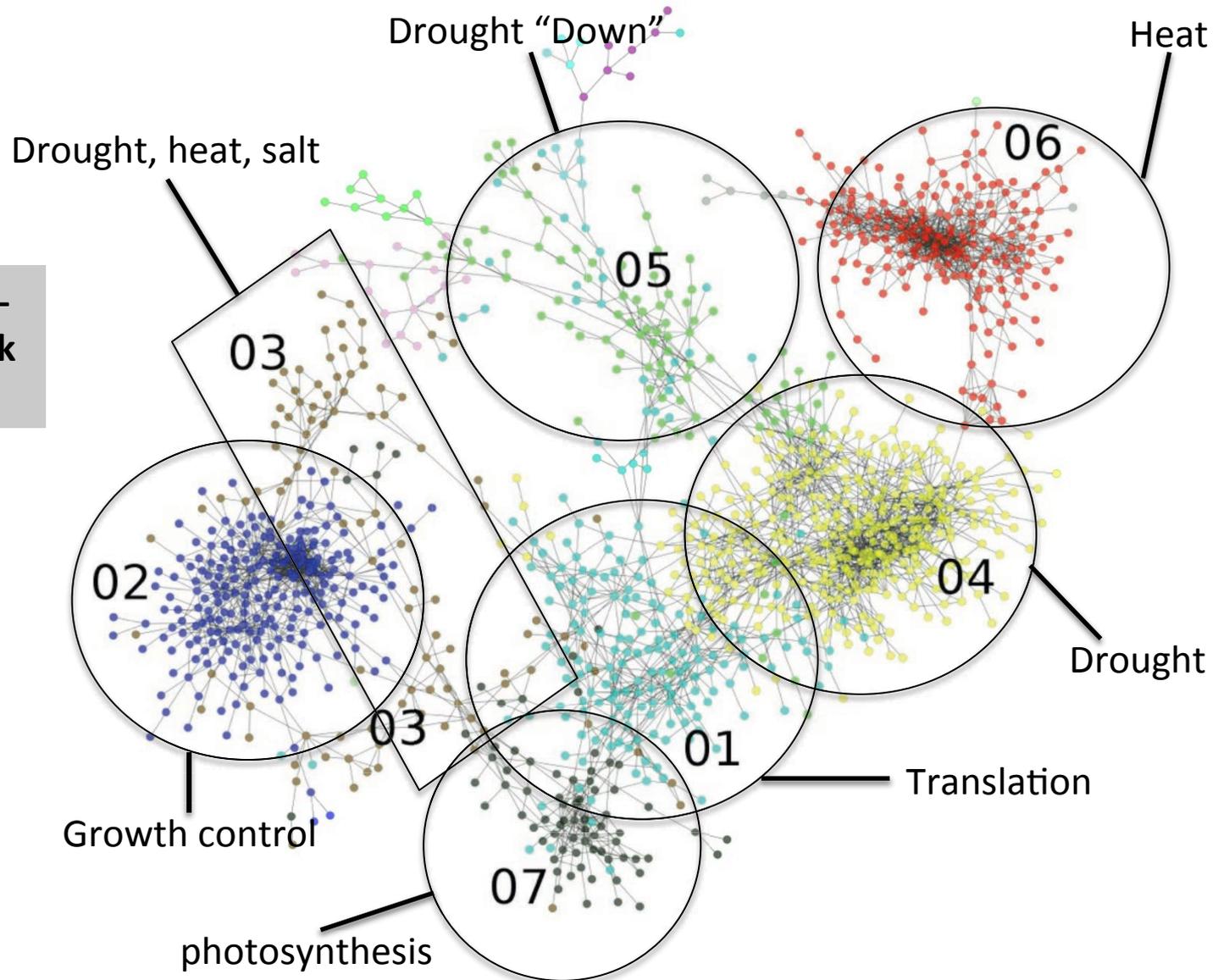


**Genome
enabled
plant
systems**

The landscape is mature with many options and developed user communities.

Example: Network analysis elucidates integration of stress responses, photosynthesis, and growth control (Brachypodium)

Weighted Gene Co-expression Network Analysis (WGCNA)



The Bottleneck Has Shifted



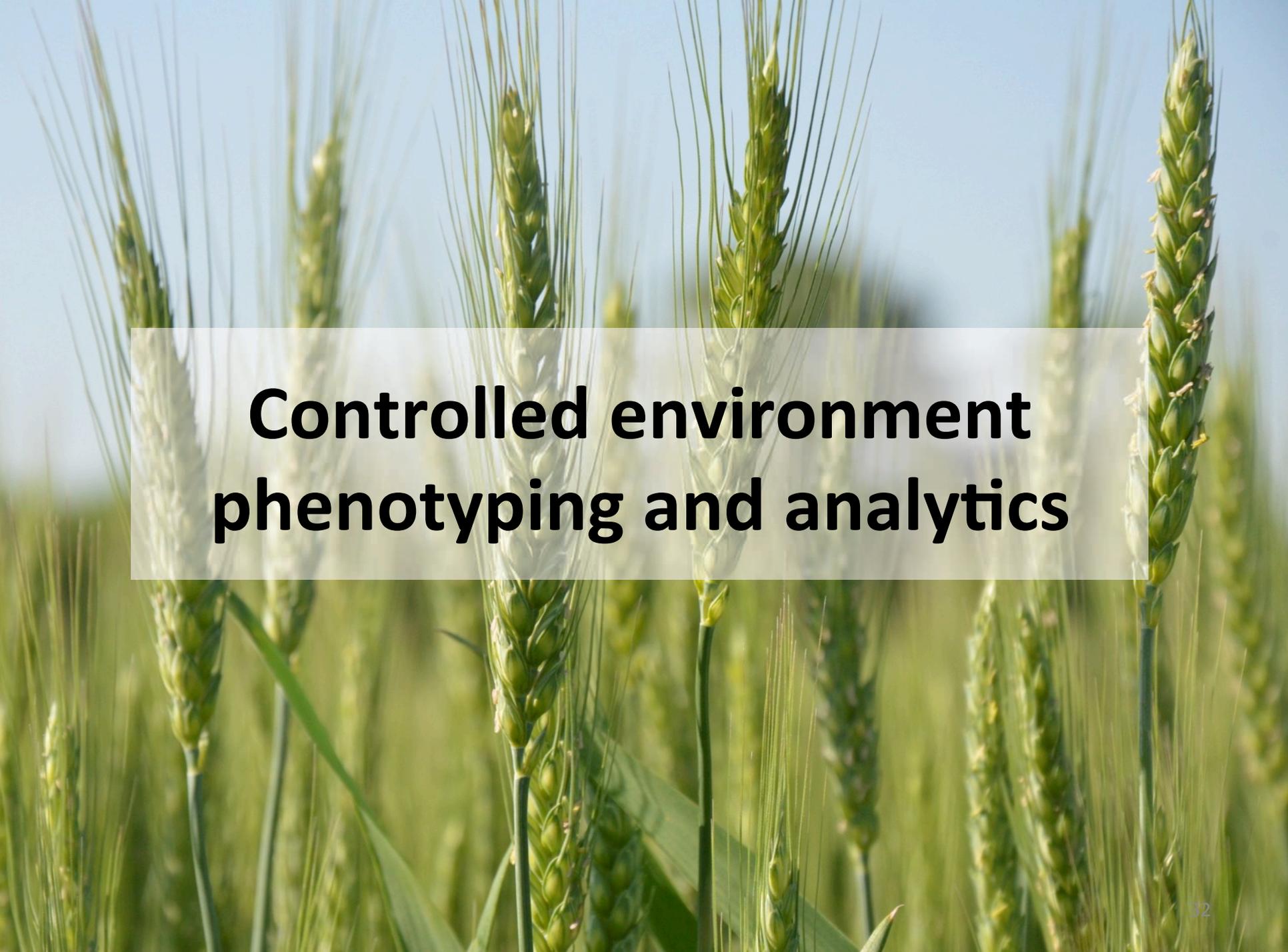
A systems-level analysis of drought and density response in the model C4 grass *Setaria viridis*

Tom Brutnell (PI); Co-PIs Andrew Leakey, Asaph Cousins, Ivan Baxter, Todd Mockler, Jose Dinneny, Sue Rhee, Dan Voytas, Hector Quemada

University of Illinois - Urbana-Champaign

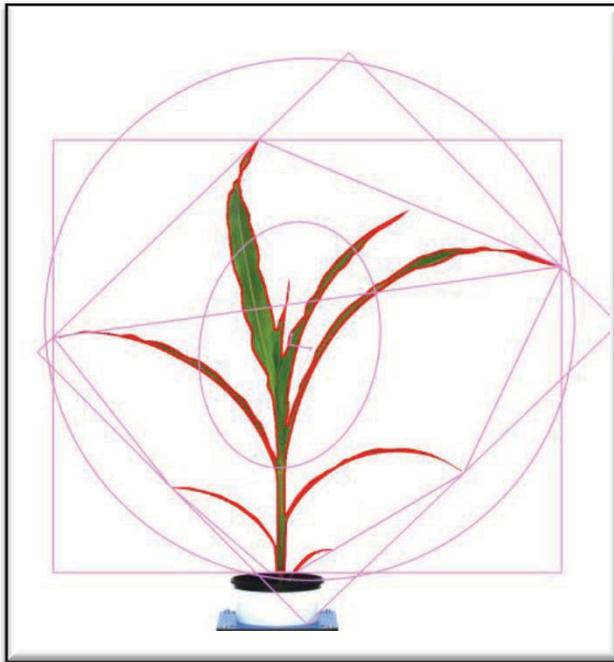


DOE Genomic Science Program - Award Number DE-SC0008769

A close-up photograph of several green wheat spikes. The spikes are in various stages of growth, with some showing more developed grains. The background is a clear, bright blue sky. The overall scene is brightly lit, suggesting a sunny day.

Controlled environment phenotyping and analytics

High-throughput image-based controlled-environment phenotyping



Visible Light Imaging

High-resolution color images for comprehensive morphological and growth phenotyping

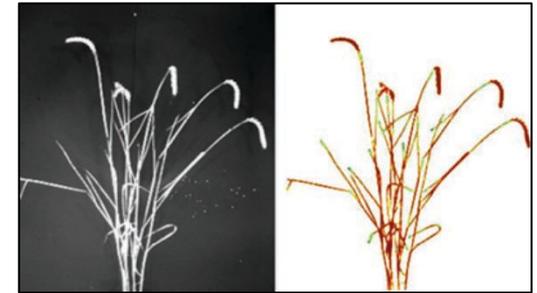


Fluorescent Light Imaging

Blue light (< 500 nm) visualize

any related fluorescence

- Chlorophyll (cont., flash)
- Green fluo. Protein (GFP)
- Phenolics
- Auto fluorescence



Near-Infrared Light (NIR) Imaging

Shoots: Measuring water distribution and dynamics

Roots: spatial distribution of water content in soil

Danforth's controlled-environment phenotyping facility

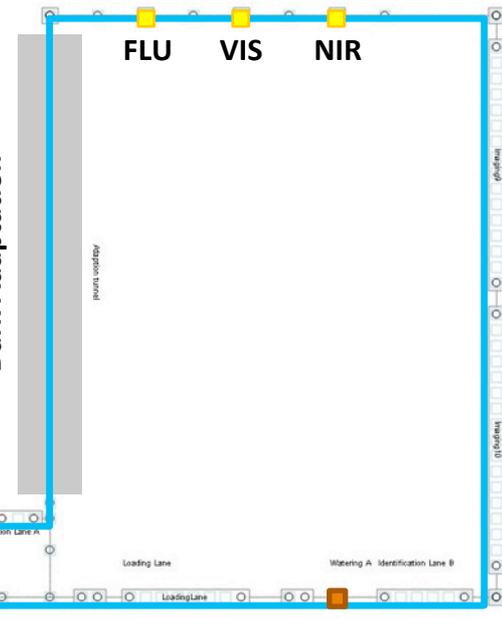
- 1200 plant capacity
- The first non-commercial high-throughput plant phenotyping facility in the US
- Imaging and software enable daily measurements of growth, 3D architecture, biomass accumulation, fluorescence and infrared signatures.



Conviron Growth House



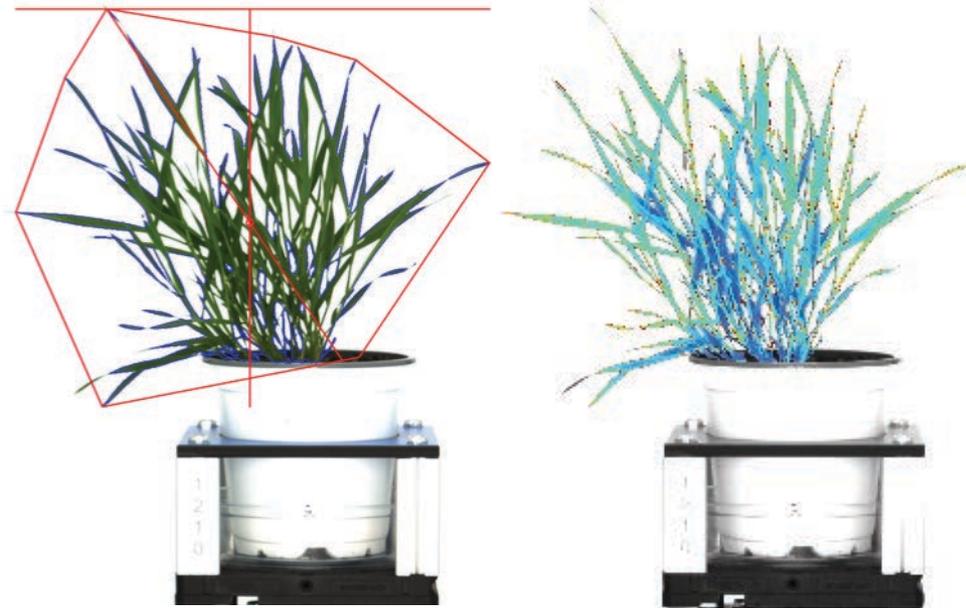
Imaging Loop



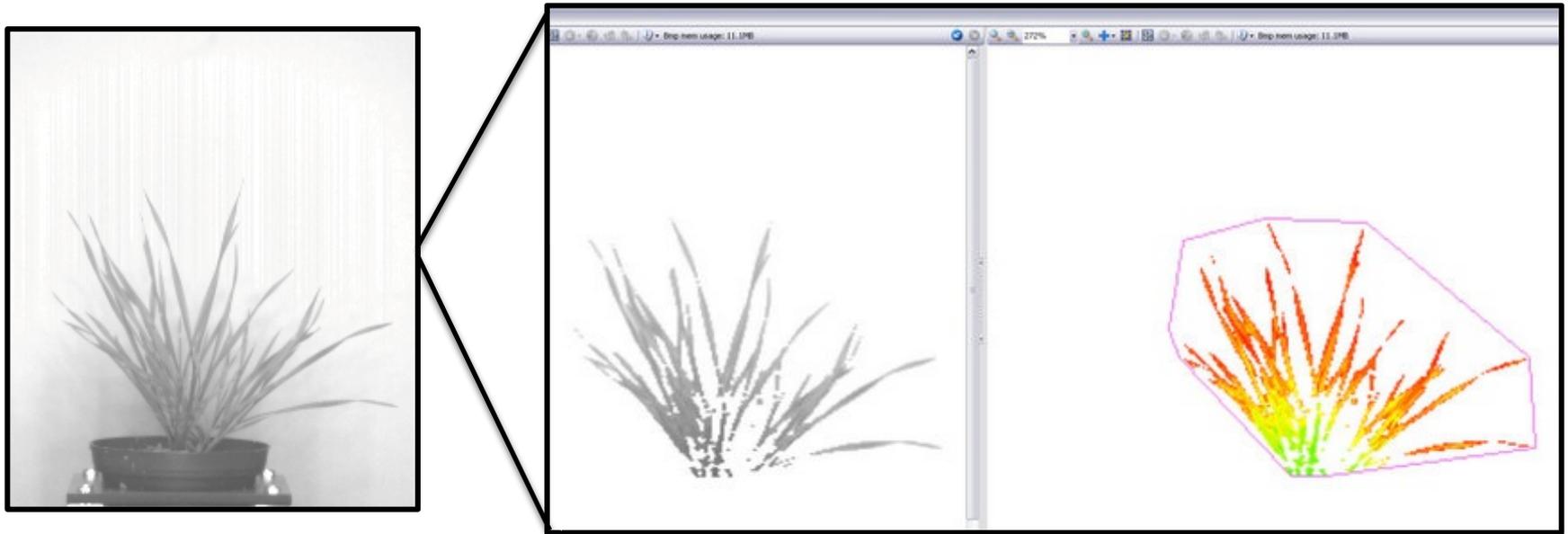
Three Imaging Stations: Fluor



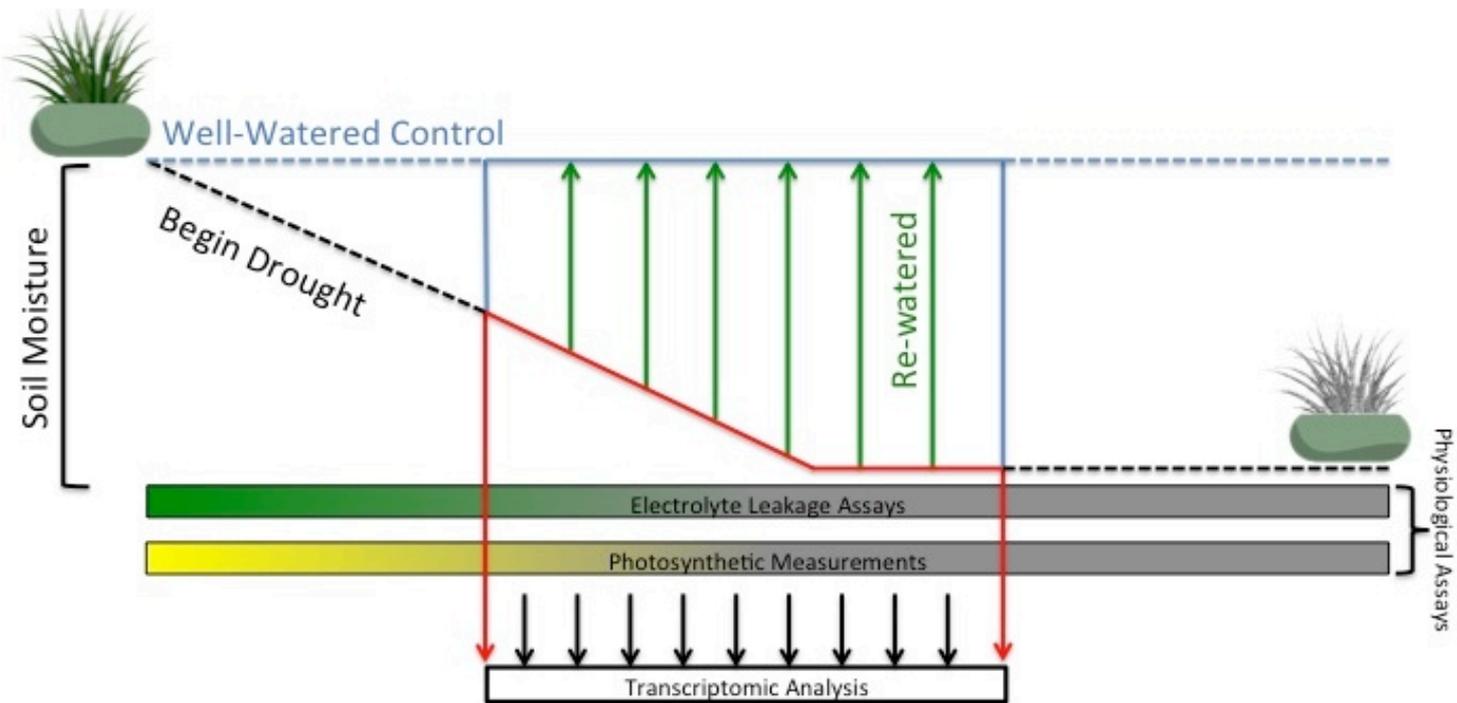
Three Imaging Stations: VIS



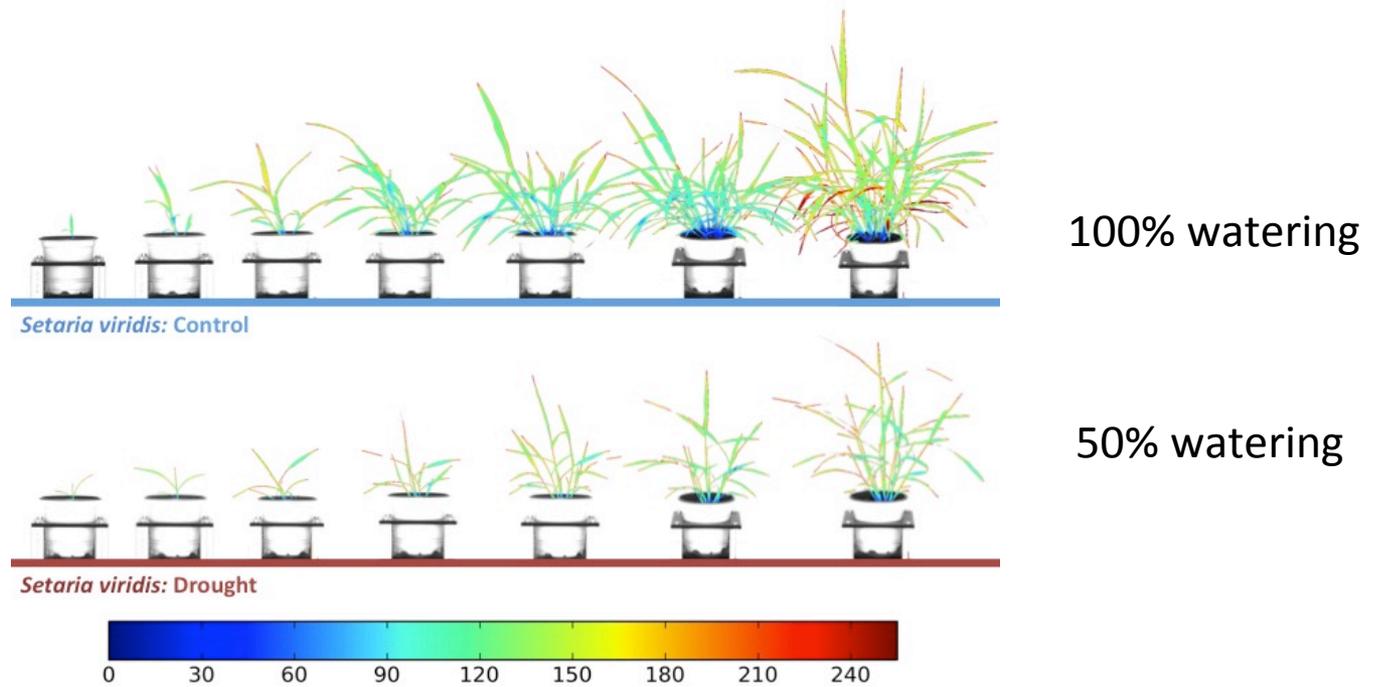
Three Imaging Stations: NIR



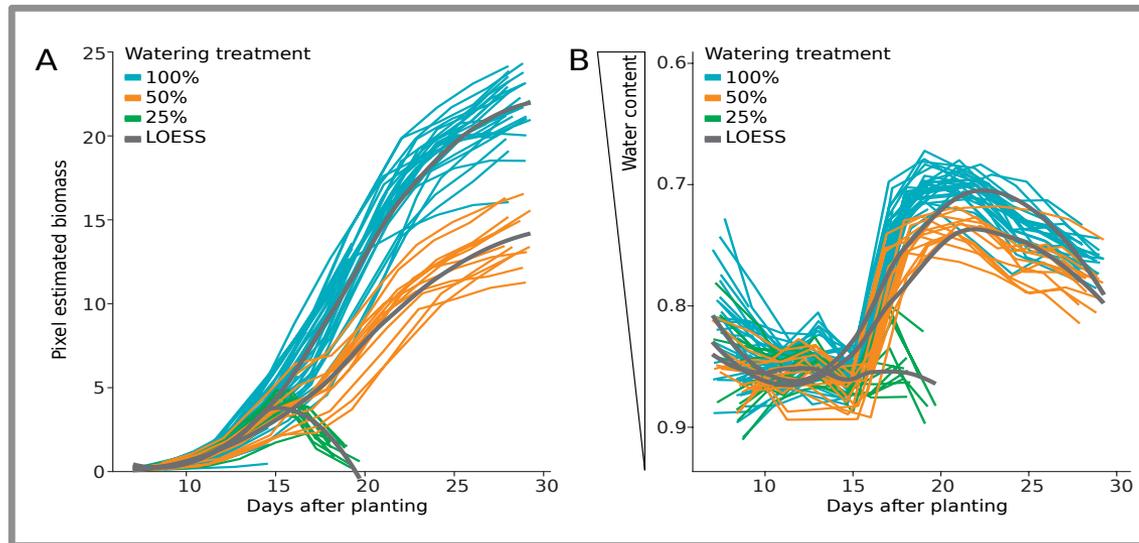
Controlled environment phenotyping enables coupling physiological analysis and 'omics profiling to understand plant-environment interactions and effects on growth and yield



Quantifying effects of drought on *Setaria viridis* growth, development, and physiology

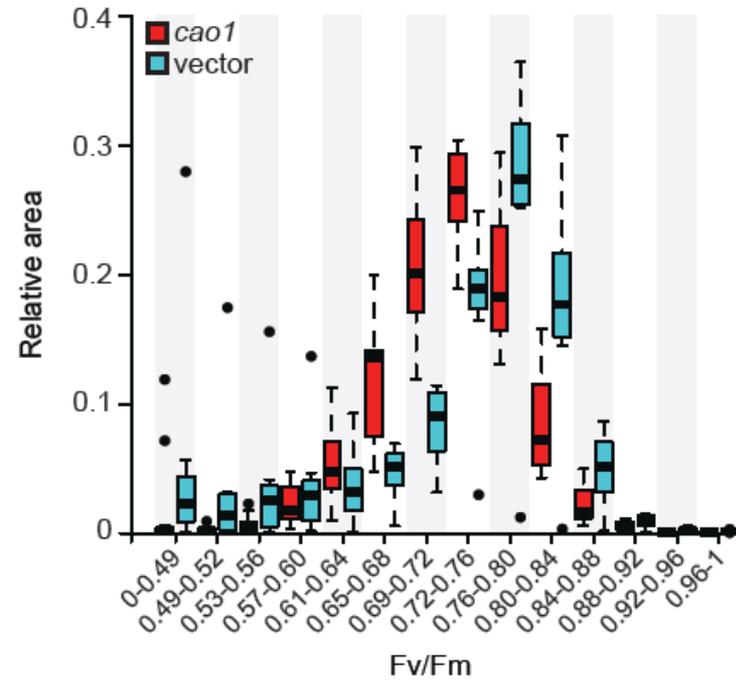
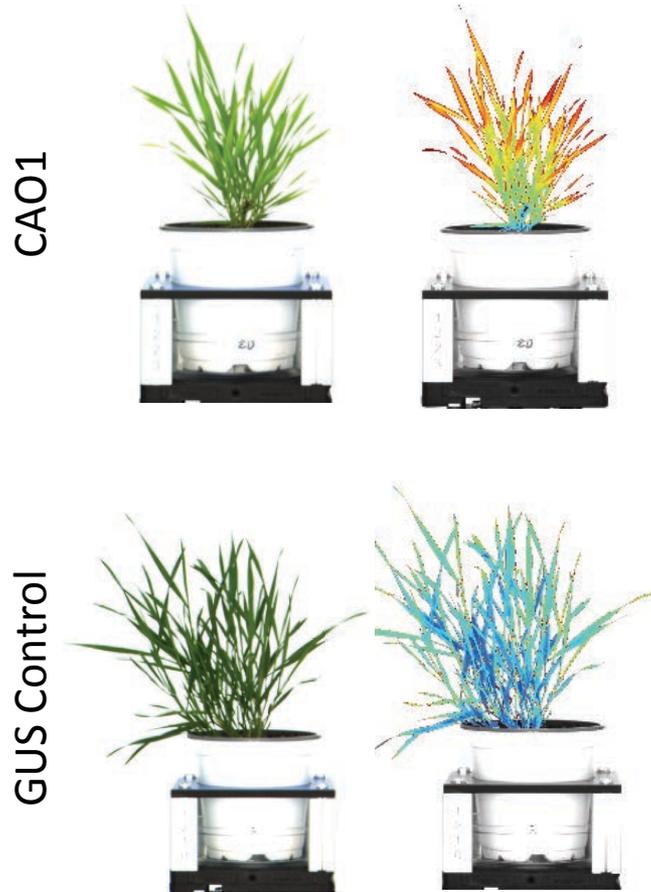


Visible
Imaging:
biomass



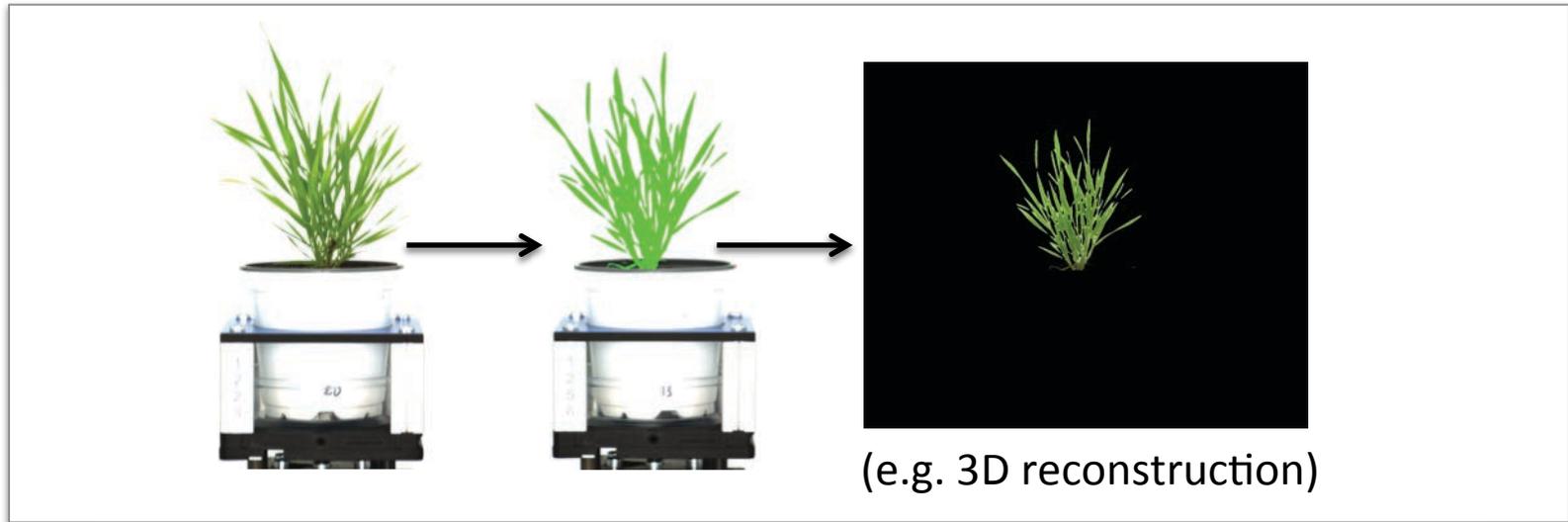
NIR
Imaging:
water
content

Fluorescence imaging-based photosynthetic analysis (benchmarking)



Lowered Fv/Fm ratio in *cao1* (chlorophyllide a oxygenase) is consistent with fewer open reaction centers available in this photosynthetic mutant.

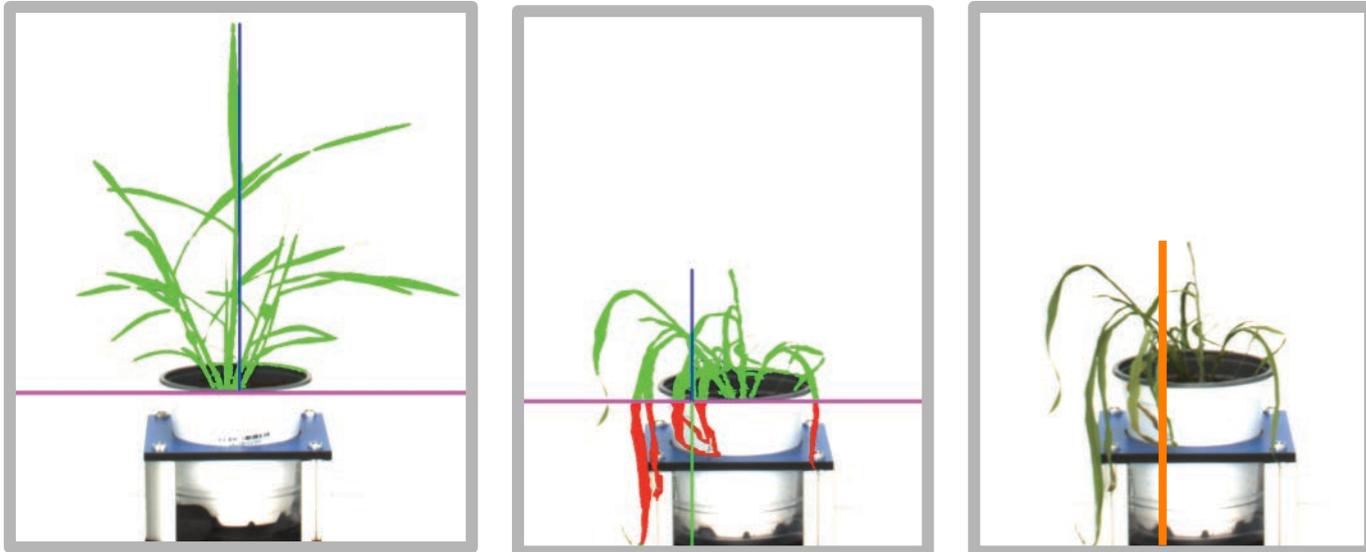
PlantCV (Plant Computer Vision)



- Built with opensource libraries OpenCV, Numpy and Matplotlib
- Integration of plant identification, feature identification, and feature analysis
- Works on other types of images, not LemnaTec-specific
- Will be published and released publicly - soon

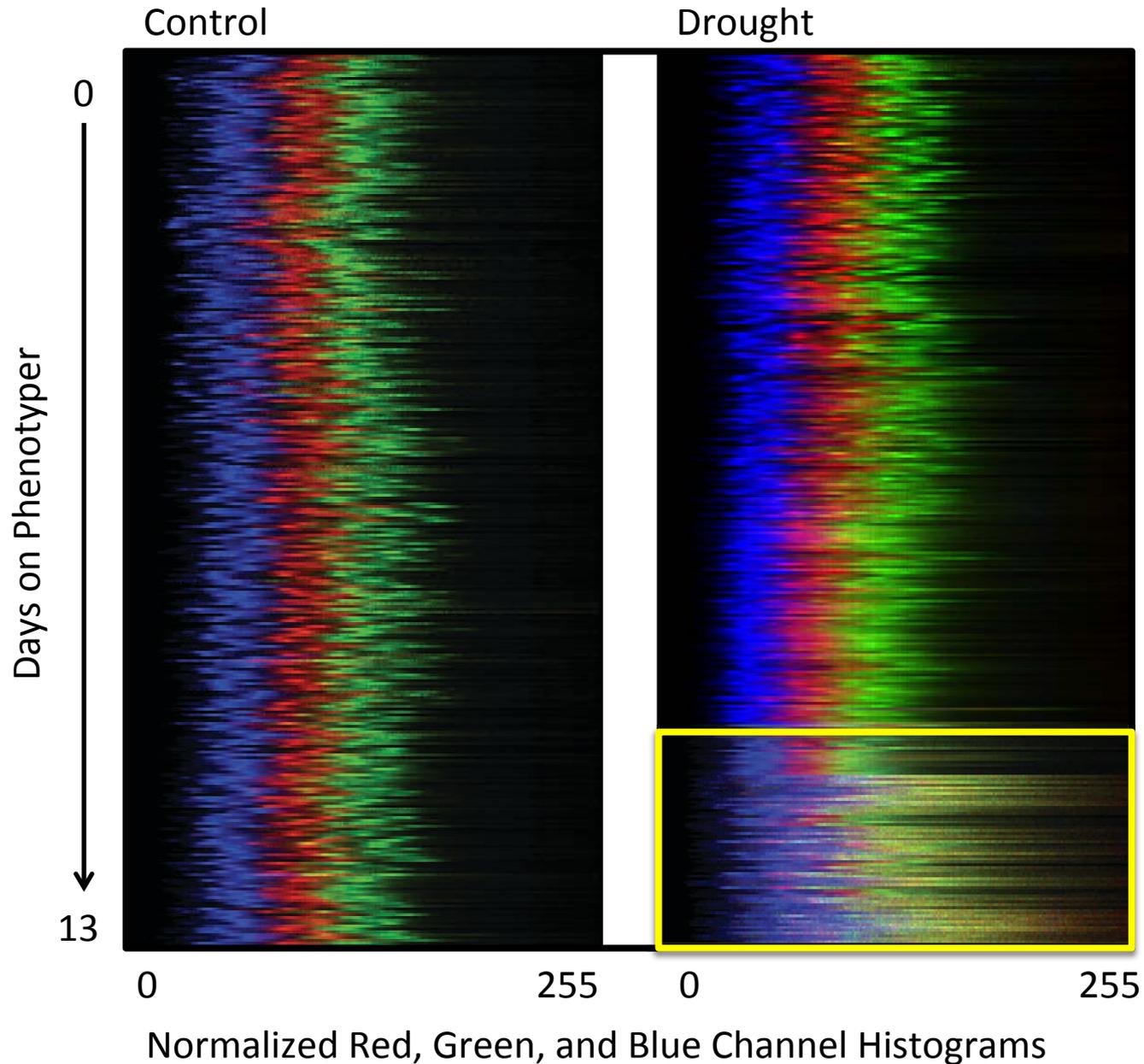
Being developed jointly by Mockler, Carrington, Baxter, and Brutnell groups at Danforth

Proxy measurements for more complex traits of interest



Automatic wilting assessment as a proxy for drought tolerance

Setaria viridis: Visualization of color spectrum histograms to develop drought signatures



A photograph of a cornfield with young, green plants in rows. The plants are growing in dark brown soil. A semi-transparent white rectangular box is overlaid in the center of the image, containing the text "Where do we need to go?".

Where do we need to go?

Where do we need to go – genomics/genetics?

- Sequencing technology and bioinformatics are well established – and not a real limiting factor now (but should be tightly integrated into improvement efforts).
- For some potential target species, existing ‘omics resources are significant but insufficient. Sequence data acquisition is still needed (diversity collections, understanding epigenomic impacts).
- For example, existing gene expression profiling resources in sorghum limit the quality and utility of gene network models that can drive hypothesis-driven allele-specific breeding strategies.
- Investment in data organization/coordination and integration is needed in specific crops – generic databases don’t address this issue for specific crops

Where do we need to go – phenomics/analytics?

- More plant phenotyping data needs to be acquired in order to inform the development of improved analytics.
- It's not all about image-based phenotyping - molecular phenotypes matter (e.g. gene expression, metabolomics).
- Algorithms and software are needed to automatically extract and/or derive digital plant phenotypes from large, typically image-based datasets.
- Hyperspectral imaging seems to have great promise – but we don't have off-the-shelf hyperspectral profile models that correlate with key plant traits.
- Robust packages/pipelines and communities of users are needed to drive standardization.
- Investment in data organization/coordination/curation and ongoing integration.